

AD-A094 604

STANFORD UNIV CA DEPT OF COMPUTER SCIENCE  
THEORY OF COMPILER SPECIFICATION AND VERIFICATION. (U)

MAY 80 W H POLAK

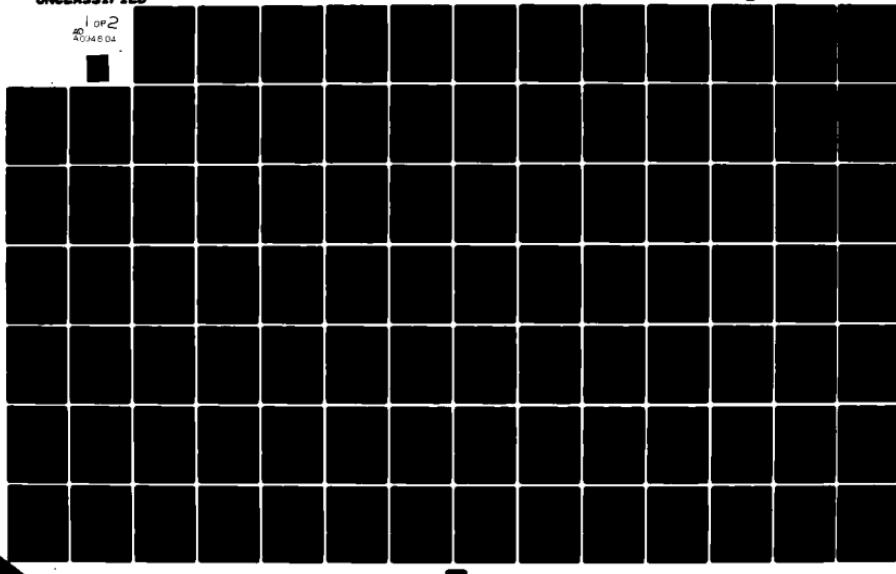
F/B 9/2

MDA903-76-C-0206

NL

UNCLASSIFIED

1 OP2  
40345 04



UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE WHEN DRAFTED

REPORT DOCUMENT		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER  PVG-17	2. GOVT ACCESSION NO.  AD-A0944684	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle)  THEORY OF COMPILER SPECIFICATION AND VERIFICATION.	5. TYPE OF REPORT & PERIOD COVERED  Thesis	
7. AUTHOR(s)  W. H. Polak	6. PERFORMING ORG. REPORT NUMBER  MDA90376C0206 MCS7600321A1	
9. PERFORMING ORGANIZATION NAME AND ADDRESS  Stanford University Stanford, California	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS  AO 2494	
11. CONTROLLING OFFICE NAME AND ADDRESS  Defense Advanced Research Projects Agency 1400 Wilson Blvd. Arlington, VA 22209	12. REPORT DATE  May 1980	
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)	13. NUMBER OF PAGES 193	
15. SECURITY CLASS. (of this report)  Unclassified		
15a. DECLASSIFICATION/DOWNGRADING SCHEDULE		
16. DISTRIBUTION STATEMENT (of this Report)  APPROVED FOR PUBLIC RELEASE. DISTRIBUTION UNLIMITED. THIS DOCUMENT MAY BE REPRODUCED FOR ANY PURPOSE OF THE U.S. GOVERNMENT		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)  NA		
18. SUPPLEMENTARY NOTES  NA		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number)  COMPUTERS COMPILER PASCAL PASCAL PLUS		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)  The formal specification, design, implementation, and verification of a compiler for a Pascal-like language is described. All components of the compilation process such as scanning, parsing, type checking, and code generation are considered.  The implemented language contains most control structures of Pascal, recursive procedures and functions, and jumps. It provides user defined data types including arrays, records, and pointers. A simple facility for input - output is provided.		

DD FORM 1 JAN 73 1473 EDITION OF 1 NOV 68 IS OBSOLETE

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

The target language assumes a stack machine including a display mechanism to handle procedure and function calls.

The compiler itself is written in Pascal Plus, a dialect of Pascal accepted by the Stanford verifier. The Stanford verifier is used to give a complete formal machine checked verification of the compiler.

One of the main problem areas considered is the formal mathematical treatment of programming languages and compilers suitable as input for automated program verification systems.

Several technical and methodological problems of mechanically verifying large software systems are considered. Some new verification techniques are developed, notably methods to reason about pointers, fixed points, and quantification. These techniques are of general importance and are not limited to compiler verification.

The result of this research demonstrates that construction of large correct programs is possible with the existing verification technology. It indicates that verification will become a useful software engineering tool in the future. Several problem areas of current verification systems are pointed out and areas for future research are outlined.

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

**Stanford Program Verification Group**  
**Report No. PVG-17**

**May 1980**

**Department of Computer Science**  
**Report No. STAN-CS-80-802**

## **THEORY OF COMPILER SPECIFICATION AND VERIFICATION**

**by**

**Wolfgang Heinz Polak**

APPROVED FOR PUBLIC RELEASE  
THIS DOCUMENT MAY BE REPRODUCED FOR  
ANY PURPOSE OF THE U. S. GOVERNMENT

**Research sponsored by**

**Advanced Research Projects Agency  
and  
National Science Foundation**

**COMPUTER SCIENCE DEPARTMENT  
Stanford University**



**81 2 03 025**

1

## THEORY OF COMPILER SPECIFICATION AND VERIFICATION

Wolfgang Heinz Polak

Department of Computer Science  
Stanford University  
Stanford, California 94305

May 1980

APPROVED FOR PUBLIC RELEASE  
THIS DOCUMENT MAY BE REPRODUCED FOR  
ANY PURPOSE OF THE U. S. GOVERNMENT

*This research was supported by the Advanced Research Projects Agency of the Department of Defense under Contract MDA903-76-C-0206, by the National Science Foundation under Contract MCS 76-00321-A1, and by the "Studienstiftung des deutschen Volkes."*

© Copyright 1980 by Wolfgang Heinz Polak

**Stanford Program Verification Group**  
**Report No. PVG-17**

**May 1980**

**Department of Computer Science**  
**Report No. STAN-CS-80-802**

## **THEORY OF COMPILER SPECIFICATION AND VERIFICATION**

**by**

**Wolfgang Heinz Polak**

### **ABSTRACT**

The formal specification, design, implementation, and verification of a compiler for a Pascal-like language is described. All components of the compilation process such as scanning, parsing, type checking, and code generation are considered.

The implemented language contains most control structures of Pascal, recursive procedures and functions, and jumps. It provides user defined data types including arrays, records, and pointers. A simple facility for input-output is provided.

The target language assumes a stack machine including a display mechanism to handle procedure and function calls.

The compiler itself is written in Pascal Plus, a dialect of Pascal accepted by the Stanford verifier. The Stanford verifier is used to give a complete formal machine checked verification of the compiler.

One of the main problem areas considered is the formal mathematical treatment of programming languages and compilers suitable as input for automated program verification systems.

Several technical and methodological problems of mechanically verifying large software systems are considered. Some new verification techniques are developed, notably methods to reason about pointers, fixed points, and quantification. These techniques are of general importance and are not limited to compiler verification.

The result of this research demonstrates that construction of large correct programs is possible with the existing verification technology. It indicates that verification will become a useful software engineering tool in the future. Several

problem areas of current verification systems are pointed out and areas for future research are outlined.

*This thesis was submitted to the Department of Computer Science and the Committee on Graduate Studies of Stanford University in partial fulfillment of the requirements for the degree of Doctor of Philosophy.*

*This research was supported by the Advanced Research Projects Agency of the Department of Defense under ARPA Order No. 2494, Contract MDA903-76-C-0206 and National Science Foundation under Contract NSF MCS 76-00321-A1. The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of Stanford University, or any agency of the U. S. Government.*

© Copyright 1980

by

Wolfgang Heinz Polak

**Abstract:**

The formal specification, design, implementation, and verification of a compiler for a Pascal-like language is described. All components of the compilation process such as scanning, parsing, type checking, and code generation are considered.

The implemented language contains most control structures of Pascal, recursive procedures and functions, and jumps. It provides user defined data types including arrays, records, and pointers. A simple facility for input - output is provided.

The target language assumes a stack machine including a display mechanism to handle procedure and function calls.

The compiler itself is written in Pascal Plus, a dialect of Pascal accepted by the Stanford verifier. The Stanford verifier is used to give a complete formal machine checked verification of the compiler.

One of the main problem areas considered is the formal mathematical treatment of programming languages and compilers suitable as input for automated program verification systems.

Several technical and methodological problems of mechanically verifying large software systems are considered. Some new verification techniques are developed, notably methods to reason about pointers, fixed points, and quantification. These techniques are of general importance and are not limited to compiler verification.

The result of this research demonstrates that construction of large correct programs is possible with the existing verification technology. It indicates that verification will become a useful software engineering tool in the future. Several problem areas of current verification systems are pointed out and areas for future research are outlined.

Accession For	NTIS CRA&I	<input checked="" type="checkbox"/>
DTIC TAB		<input type="checkbox"/>
Unclassified		<input type="checkbox"/>
Justification		
By		
Distribution		
Availability Codes		
Air Mail		
Dist	Special	

## Preface

About three years ago David Luckham hinted to me the possibility of verifying a "real" compiler. At that time the idea seemed unrealistic, even absurd. After looking closer at the problem and getting more familiar with the possibilities of the Stanford verifier a verified compiler appeared not so impossible after all. In fact, I was fascinated by the prospect of creating a large, correct piece of software; so this subject became my thesis topic. I am very grateful to David Luckham for suggesting this topic and for his continued advice.

The research has drastically changed my view of verification and programming in general. Analysis and design of programs (even large ones) can be subject to rigorous mathematical treatment - the art of programming may become a science after all.

Naturally, the reader will be skeptical, still, I hope to be able to convey some of my fascination.

This work would have been impossible without the use of the Stanford verifier. I have to thank all members of the Stanford verification group for providing this excellent tool. Don Knuth's text processing system TEX was a most valuable asset for typesetting a manuscript that must be any typist's nightmare.

I thank my reading committee David Luckham, Zohar Manna, and Susan Owicki for their valuable time, for their careful reading, and their helpful advise. Friedrich von Henke contributed through numerous discussions and careful perusal of initial drafts of my writing. Last but not least I thank my wife Gudrun for her support.

<b>Chapter I. Introduction</b>	
1. Overview	1
2. Program verification	2
2.1. Writing correct programs	3
2.1.1. Program design	4
2.1.2. An example	5
2.2. Logics of programming	6
2.2.3. Logic of computable functions (LCF)	7
2.3.0. First order logic	7
2.3.0. Hoare's logic	8
2.3.0. Comparison	8
2.3. The Stanford verifier	8
3. Formal definition of programming languages	10
3.1. Syntax	11
3.1.1. Micro syntax	11
3.1.2. Phrase structure	12
3.2.1. Tree Transformations	13
3.2. Semantics	13
3.2.3. Operational semantics	13
3.2.3. Denotational semantics	14
3.2.4. Floyd - Hoare logic	14
3.2.4. Algebraic semantics	15
3.2.5. Others	16
4.0. Machine languages	17
4.0. Summary	17
4. Developing a verified compiler	17
4.1. What we prove	18
4.2. Representation	19
4.4. Scanner	20
4.4. Parser	21
4.5. Semantic analysis	22
4.6. Code generation	23
4.6.1. Necessary theory	23
4.6.2. Implementation	24
5. Related work	25
5.1. Previous work on compiler verification	25
5.3. Relation to our work	26
5.3. Compiler generators	27
6. Organization of this thesis	28
<b>Chapter II. Theoretical framework</b>	
1. Basic concepts	30
1.1. Functions	30
1.2. First order logic	31
1.2.2. Syntax	31
1.2.2. Semantics	31
2. Scott's logic of computable functions	32
2.2. Basic definitions	32
2.2. Operations on domains	33
2.3. Conditionals	35
3.1. Lists	36
3. Induction proofs	36
3.1. Fixed point induction	36
3.2. Recursion induction, structural induction	37
4. Verification techniques	38
4.1. Pointers	38
4.1.1. Pointers in the Stanford verifier	38
4.1.2. Reasoning about pointers	39
4.1.3. Reasoning about extensions	40
4.2. Quantification	41
4.3. Computable functions and first order logic	43
4.3.1. Standard interpretations	43
4.3.2. Higher order functions	44
4.3.3. Least fixed points	45
5. Representations	46
5.1. Primitive types	46

5.3. Complex types	47	5.1. Syntax of <i>LS</i>	67
6.1. Recursive domains	47		
<b>6. Fixed points</b>	<b>48</b>	<b>5. Tree transformations</b>	<b>67</b>
6.1. Reasoning about fixed points	48	5.1. Abstract syntax	67
6.2. Operationalisation of fixed points	49	5.1.1. Syntactic domains	67
<b>7. Intermediate representation of programs</b>	<b>51</b>	6.0. Tree transformations for <i>LS</i>	69
7.2. Trees	51	<b>6. Semantics of <i>LS</i></b>	<b>69</b>
7.2.0. $\Sigma$ -trees on <i>X</i>	51	6.1. Semantic concepts	70
7.2.0. Operations on trees	52	6.1.1. Semantic domains	70
7.2. Abstract syntax	52	6.1.2. Types and modes	74
7.2.2. Constructor and selector functions	53	6.1.3. Auxiliary functions, static semantics	75
7.2.2. Conformity relations	53	6.2. Static semantics	76
		6.2.1. Declarations	76
		6.2.3. Types	77
		6.2.5. Labels and identifiers	77
		6.2.5. Expressions	78
		6.2.5. Statements	78
		6.3.1. Commands	79
		6.3. Dynamic semantics	79
<b>Chapter III. Source and target languages</b>		6.3.1. Auxiliary functions	79
1. The source language	55	6.3.2. Memory allocation	80
1.1. Data structures	55	6.3.3. Declarations	81
1.2. Program structures	57	6.3.4. Expressions	82
<b>2. Formal definition of <i>LS</i></b>	<b>58</b>	6.3.5. Statements	83
2.1. Structure of the formal definition	58		
2.2. Denotational semantics	59		
2.2.2. General concepts	59		
2.2.2. Semantic concepts of Algol-like languages	60		
2.2.4. Denotational definition of a machine language	62		
2.2.4. Notational issues	62		
<b>3. Micro syntax</b>	<b>63</b>	<b>7. The target language <i>LT</i></b>	<b>84</b>
3.1. Definitional formalism	63	7.1. A hypothetical machine	85
4.1. Micro syntax of <i>LS</i>	64	7.1.2. Design decisions	85
<b>4. Syntax</b>	<b>65</b>	7.1.2. Architecture	85
4.1. Labeled context free grammars	65	7.1.3. Instructions	87
4.1.3. The accepted language	66	7.2. Formal definition of <i>LT</i>	88
4.1.3. Parse trees	66	7.2.2. Abstract syntax	88
4.1.3. The function defined by a labeled grammar	66	7.2.2. Semantic domains	88
		7.2.3. Semantic equations	89

#### **Chapter IV. The compiler proof**

1. Verifying a compiler . . . . .	90	5. Code generation . . . . .	122
1.1. The compiler . . . . .	90	5.1. Principle of code generation . . . . .	123
1.1.2. Correctness statement . . . . .	90	5.2. Modified semantics definitions . . . . .	125
1.1.2. Structure of the compiler . . . . .	90	5.2.1. A structured target language . . . . .	125
1.2. The individual proofs . . . . .	92	5.2.2. A modified definition of $LS$ . . . . .	127
2. A scanner for $LS$ . . . . .	93	5.3. Relation between $LS$ and $LT$ . . . . .	129
2.1. Underlying theory . . . . .	93	5.3.2. Compile time environments . . . . .	130
2.1.1. A suitable definition . . . . .	93	5.3.2. Storage allocation . . . . .	130
2.1.2. Axiomatisation of concepts . . . . .	94	5.3.3. Storage maps . . . . .	131
2.2. Basic algorithm . . . . .	95	5.3.4. Relations between domains . . . . .	132
2.3. Implementation details . . . . .	97	5.3.5. Existence of recursive predicates . . . . .	135
3. A parser for $LS$ . . . . .	98	5.4. Implementation of the code generation . . . . .	136
3.1. LR theory . . . . .	98	5.4.1. Specifying code generating procedures . . . . .	136
3.1.1. LR-parsing tables . . . . .	99	5.4.2. Treatment of labels . . . . .	140
3.1.2. The LR-parsing algorithm . . . . .	101	5.4.4. Declarations . . . . .	143
3.2.0. Axiomatisation . . . . .	102	5.4.4. Procedures and functions . . . . .	144
3.2. Tree transformations . . . . .	102	5.4.6. Blocks . . . . .	146
3.3.1. Building abstract syntax trees . . . . .	104	5.4.6. Refinement . . . . .	146
3.3. Refinement . . . . .	104		
3.3.2. Development cycle . . . . .	104		
3.3.2. Representation . . . . .	105		
3.3.3. Reference classes and pointer operations . . . . .	106		
4. Static semantics . . . . .	108	1. Summary . . . . .	148
4.1. Recursive declarations . . . . .	108	2. Extensions . . . . .	149
4.1.1. Operationalization . . . . .	108	2.1. Optimisation . . . . .	149
4.1.2. Revised definition of $t$ and $dt$ . . . . .	110	2.3. Register machines . . . . .	149
4.1.3. Representation of $U_s \rightarrow U_s$ . . . . .	112	2.3. New language features . . . . .	150
4.1.4. Resolving undefined references . . . . .	113	2.4. A stronger correctness statement . . . . .	151
4.2. Development of the program . . . . .	114	3. Future research . . . . .	152
4.2.1. Computing recursive functions . . . . .	114	3.2. Structuring a compiler . . . . .	152
4.2.3. Refinement . . . . .	117	3.2. Improvements of verification systems . . . . .	153
4.2.3. Representation . . . . .	118	3.4. Better verification techniques . . . . .	155
4.2.4. Auxiliary functions . . . . .	121	3.4. Program development systems . . . . .	155
5.0.0. The complete program . . . . .	122		
			References
			157

**Appendix 1. Formal definition of  $LS$** 

1. Micro Syntax of $LS$ . . . . .	167
1.3. Domains . . . . .	167
1.3.1. Languages $L_i$ . . . . .	167
1.3. Auxiliary definitions . . . . .	167
1.4. Semantic Functions . . . . .	169
2. Syntax of $LS$ . . . . .	170
3. Abstract syntax . . . . .	173
3.1. Syntactic Domains . . . . .	173
3.2. Constructor functions . . . . .	174
4. Tree transformations . . . . .	175
4.2. Programs . . . . .	175
4.2. Declarations . . . . .	176
4.3. Expressions . . . . .	177
5.1. Statements . . . . .	178
5. Semantics of $LS$ . . . . .	178
5.2. Semantic Domains . . . . .	178
6.1. Types and Modes . . . . .	179
6. Static Semantics of $LS$ . . . . .	180
6.1. Auxiliary functions, static Semantics . . . . .	180
6.2. Declarations . . . . .	185
6.4. Expressions . . . . .	187
6.4. Statements . . . . .	188
7. Dynamic Semantics of $LS$ . . . . .	190
7.1. Auxiliary functions . . . . .	190
7.2. Declarations . . . . .	196
7.3. Expressions . . . . .	199
7.4. Statements . . . . .	200

3. Semantic Domains . . . . .	203
3. Auxiliary Functions . . . . .	203
4. Semantic Equations . . . . .	204

**Appendix 3. The Scanner**

1. Logical basis . . . . .	207
1.1. Definition of the micro syntax . . . . .	207
1.2. Representation functions . . . . .	209
1.3. Sequences . . . . .	210
2. The program . . . . .	211
3. Typical verification conditions . . . . .	219

**Appendix 4. The Parser**

1. Logical basis . . . . .	223
1.2. Representation functions . . . . .	223
1.2. LR theory . . . . .	223
1.3. Tree transformations . . . . .	224
2.0. Extension operations . . . . .	225
2. The program . . . . .	225

**Appendix 5. Static semantics**

1. Logical basis . . . . .	237
1.1. Rules for $\epsilon$ . . . . .	237
1.2. Recursive types . . . . .	238

1.2.1. Types	239
2.0.0. Fixed points	241
2. The program . . . . .	241
2.1. Declarations	242
2.1.1. Types	242
2.1.2. Abstract syntax	243
2.1.3. Auxiliary functions	246
2.2. Expressions	249
2.3. Types	252

#### **Appendix 6. Code generation**

1. Logical basis . . . . .	257
2. The program . . . . .	262
2.2. Declarations	262
2.2. Virtual procedures	262
2.3. Auxiliary functions	264
2.4. Abstract syntax, types and modes	265
2.5. Code generating functions	266
2.6. Expressions	273
2.7. Commands	277
2.8. Statements	279

**Chapter I. Introduction**

*"The ultimate goals (somewhat utopian) include error-free compilers, . . ."*

S. Greibach

**1. Overview**

In this thesis we describe the design, implementation, and verification of a compiler for a Pascal-like language. While previous work on compiler verification has focused largely on proofs of abstract "code generating algorithms" we are interested in a "real" compiler translating a string of characters into machine code efficiently. We consider all components of such a compiler including scanning, parsing, type checking and code generation.

Our interest is twofold. First, we are concerned with the formal mathematical treatment of programming languages and compilers and the development of formal definitions suitable as input for automated program verification systems. Second, we are interested in the more general problem of verifying large software systems.

There are several reasons for verifying a compiler. Compilers are among the most frequently used programs. Consequently, if we invest in program proofs it is reasonable to do this in an area where we may expect the highest payoff. Verification techniques are in general applied to programs written in high level languages. If we ever want to use verified programs we have to be able to correctly translate them into executable machine languages; another reason why correct compilers are a necessity.

The implemented language, *L<sub>T</sub>*, contains all features of Pascal [JW76] that are of interest to compiler constructors. The language contains most control structures of Pascal, recursive procedures and functions, and jumps. It provides user defined data types including arrays, records, and pointers. A simple facility for input - output is included; each program operates with one input and one output file.

The target language, *L<sub>T</sub>*, assumes a stack machine including a display mechanism [RR64, Or73] to handle procedure and function calls. This language simplifies our task somewhat as it avoids the issue of register allocation and similar irrelevant details. But at the same time the target language is realistic in that similar machine architectures exist, notably the B600.

The compiler itself is written in Pascal Plus, a dialect of Pascal accepted by the Stanford verifier.<sup>1</sup> The Stanford verifier [SV79] is used to give a complete formal machine checked verification of the compiler.

We review existing methods for the formal definition of programming languages. We use the most appropriate definitional methods to formally define source and target language.

We further investigate how the correctness of a compiler can be specified in a form suitable for mechanical verification. Here, we have to deal with several technical issues such as fixed points and reasoning about pointers.

We show that an efficient program can be developed systematically from such specifications.

During this research verification has proven to be a most useful tool to support program development; verification should be an integral part of the development process and plays a role comparable to that of type checking in strongly typed languages. This methodology of verification supported programming is not limited to compilers, rather it is applicable to arbitrary problem domains.

The results of this thesis are encouraging and let us hope that verification will soon become a widely accepted software engineering tool. But also this work reveals many of the trouble spots still existing in today's verification technology. We point to several promising research areas that will make verification more accessible. The need for better human engineering and integrated software development systems is particularly urgent.

**2. Program verification**

Verification has provoked several controversial statements by opponents and proponents recently [DL79]. Therefore it is appropriate at this point to clarify what verification is, what it can do, and, most importantly, what it cannot do.

1.) Notable difference to standard Pascal is that formal documentation is an integral part of the language, for more details see 2.3.

To verify a program means to prove that the program is consistent with its specifications. Subsequently, the term "verification" is used as a technical term referring to the process of proving consistency with specifications. A program is "verified" if a consistency proof has been established.

Formal specifications for a program can express different requirements. For example, a specification can be "the program terminates for each set of input data". Or, even more trivially, one can specify that "the program satisfies all type and declaration constraints"; every compiler verifies this property. But of course, we consider more interesting properties of our compiler. What exactly its specifications are is discussed in the following sections and in more detail in chapter IV.

Since specifications can be weak and need not (and generally do not) express all requirements for a program, verification should not be confused with correctness in the intuitive sense (i.e., the program does what one expects it to do). Verification is not a substitute for other software engineering techniques such as systematic program development, testing, walk-through and so on; rather verification augments these techniques. It gives us an additional level of confidence that our program has the property expressed in its specifications.

Depending on the application of a program certain errors may be mere annoyances while other may have disastrous effects. We can classify errors as "cheap" and "expensive". For example, in terms of our compiler cheap properties are the reliability of error recovery and termination. An expensive error would be if the compiler would translate an input program without reporting an error but would generate wrong code.

As long as verification is expensive we can concentrate our efforts on the most costly requirements of a program. These requirements can be formalized and taken as specifications for the program. Verification can be used to guarantee these expensive requirements. Other cheaper properties can be validated by conventional testing methods. Redundant code can easily be added to a verified program to increase its reliability or establish additional properties without affecting verified parts of the program. For example, in our compiler a sophisticated error recovery could be added without invalidating program proofs.

## 2.1. Writing correct programs

The non-technical use of the expression "to verify a program" suggests that there is an existing program which we subject to a verification. This is possible in principle but most certainly not practical for large software systems;

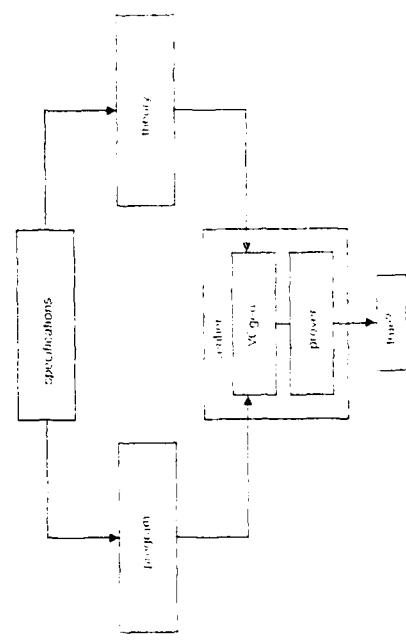


Fig. 1

it will only work for "toy" programs.

Instead, verification should be an integral part of program development: a program and its proof should be developed simultaneously from their specifications.

### 2.1.1. Program design

The process of designing a verified program can be visualized by figure 1. The starting point for the program design is a formal specification of the problem. In general, this specification alone is insufficient to prove the correctness of an implementation. We need additional theorems and lemmas about properties of the problem domain which enable an efficient implementation. We refer to these theorems and lemmas together with the formal specifications as "theory" or logical basis for our program proof.

The original specifications as well as derived lemmas and theorems are initially given in some suitable mathematical language. They have to be encoded in a form accepted by the verification system used. This step may not be as trivial as it appears on first sight. Problems arise if the "assertion language" used by the verification system is less expressive than the mathematical language. In the case of the Stanford verifier the assertion language is quantifier

free first order logic. The compiler on the other hand is specified using Scott's logic of computable functions. In chapter II we deal with the problem of encoding the language definition in our assertion language, of particular interest is the treatment of the least fixed point operator.

Given a suitable theory we implement a program in the style of structured programming [Dv76]. Starting with a rough outline of the program we successively refine it until we derive an executable version. We start out with a first draft of our program where we are merely interested in the overall structure; all implementation details are omitted. The initial draft is then checked for consistency by the verification system. If consistency is not given, then either the program has to be corrected or stronger lemmas have to be added to the logical basis. After a successful verification we are sure that our initial design is correct (with respect to its specifications). Observe that this situation differs crucially from that in ordinary top down programming. In the latter case we have no means of validating an initial design until the refinement of the program is completed down to the lowest level and the program can actually be put on a machine and tested.

Given a verified first draft, we refine this program by implementing some of the points still left open. This step may require additional theorems and lemmas to be proven to enable efficient implementation and follows along the same lines as the initial design. We repeat this process until we arrive at an executable program.

### 2.1.2. An example

Let us explain the importance of the theory with a trivial gcd example. Our specification might be to implement a function  $f(x, y)$  such that  $f(x, y) = \text{gcd}(x, y)$ . If we are given the following formal definition of gcd an efficient implementation of  $f$  is not easily derived:

$$\text{gcd}(x, y) = \max\{z \mid z \text{ mod } z = 0 \wedge y \text{ mod } z = 0\}$$

The straightforward implementation suggested by this definition is to compute the set given in the definition and then search for the maximum. A much more efficient program is possible if we can prove additional theorems; for example  $\text{gcd}(x, y) = \text{gcd}(x \text{ mod } y, y)$ .

In general, any efficient program is based on a mathematical property of the problem domain. The really creative part of programming is to uncover relevant properties of the specifications that lead to good implementations.

In many cases we can resort to well known results in literature. But note, that if we prove a theorem, then this is a purely mathematical activity and has

nothing to do with mechanical program verification. For example, in the case of the gcd above relevant theorems are part of number theory. In the future it is to hope that we have proof checking systems which can aid us in this part of program development.

People unfamiliar with verification often do not realize the importance of theory in writing programs. Programming is misleading in this respect because it allows us to write programs without having completely understood the problem. Verification makes us aware of all the mathematical properties (or alleged properties) we use in a program. A frequent complaint about verifications is that even simple problems require a discouraging amount of mathematics. We feel that this is inherent to programming; if we want correct programs we have to understand the problem and ultimately only a mathematical treatment can help us in this.

Also we should note that "toy" verification examples are often mathematically very subtle;<sup>2</sup> many programs written in practice are based on much simpler foundations.

### 2.1.3. A word of warning

A verified program may fail to execute safely for several reasons. Firstly, a program proof could be wrong. There are two main causes for this. It is most likely that an error occurred in manually proving theorems and lemmas that are the basis of the implementation. It can be hoped, that in the future these steps are amenable to automation or can at least be checked with proof checking systems. But then of course, there is the verification system itself which could be wrong. People have pointed out, that trying to verify the verifier is a vicious circle and cannot work. True, but, we have a much simpler way out. We can take our program with formal specifications and subject it to a different verification system. Unless both verifiers contain the same bug this will reveal any error.

Secondly, even if the program proof were valid the program will ultimately be executed in a real world environment with all its uncertainties. Still, verification is an extremely promising tool to fight today's "software crisis". It is a tool to increase our confidence in a program to an extent not achievable with any other known means.

## 2.2. Logics of programming

So far we have discussed how to use verification during program development.

2.) For example dealing with permutations [Po79]

ment. We now look into the details of automatic verification systems. First, we give a general overview over existing formal systems that allow mechanical verification and discuss their advantages and disadvantages. We then take a closer look at the Stanford verifier and describe its use.

The focus of our attention is on three methods for verifying programs. These methods are selected because they appear to be reasonably well developed and more importantly, there are implementations of verification systems based on these methods.

### 2.2.1. Logic of computable functions (LCF)

LCF [Mi72a] is a formal system based on Scott's logic of computable functions.<sup>3</sup> Since in LCF programs are objects of the formal language it is possible to express very general theorems about programs in LCF. Examples are program equivalence, program transformations as well as consistency and termination.

Although LCF is best suited to reason about denotational language definitions [AA74, AA77, He75] it is also capable of expressing Hoare style proof rules (see below).

Semi-automatic proof checking systems for LCF have been implemented [Mi72b, GM75] and used for several program proofs [MW72, Co79a, Co79b, Ne73].

### 2.2.2. First order logic

In [Ca76, CM79] McCarthy and Cartwright propose to consider recursive programs as objects in a first order language. The logical system has to be augmented with an undefined symbol to account for nonterminating computations; recursive function definitions are translated into equations in the logic. Suitable induction schemas allow reasoning about least fixed points. The system has the drawback that it is restricted to first order functions. Furthermore punning between logic and programming language severely restricts the languages that can be handled and at this point the method seems not suitable to prove properties of Algol-like programs. FOL [WT74] is a proof checking program that allows mechanized reasoning about first order logic. FOL has been used for several program proofs based on McCarthy's and Cartwright's method.

<sup>3.)</sup> Details of Scott's logic are discussed in chapter II.

### 2.2.3. Hoare's logic

In [Ho69] Hoare devised a deductive system<sup>4</sup> in which programs together with specifications can be derived. The verification problem is to determine if there exists a derivation of the given program together with its specifications. Given certain intermediate documentation (invariants) there is an effective way to construct a set of first order formulas which entail that such a derivation exists. Thus, if all these formulas ('verification conditions') can be proven, the program is consistent with its specifications.

Hoare's logic is limited to partial correctness proofs, termination can be expressed by modifying the input program to include counters for all loops and suitable assertions as is outlined in [LS75, MP73]. Here, too, is a pun in that one considers expressions in the programming language to be terms in the underlying logic. This effectively excludes certain language constructs, e.g. expressions that cause side effects. The method requires the programmer not only to give input and output specifications but also specifications at intermediate points in the program.

### 2.2.4. Comparison

All of the methods mentioned above are interesting. See for example [HL74, Do76, Po80].

Verification in this thesis is exclusively based on Hoare's logic and its implementation in the Stanford Verifier. This choice has been made because the method of expressing induction hypotheses as invariants in the program is natural and powerful. Furthermore the theorem proving part of the Stanford system is extremely sophisticated which is crucial in view of the size of the program we set out to verify.

The main disadvantage of using the Stanford verifier is that proofs of theorems from specifications (as outlined above) has to be done manually. In other systems, e.g. LCF and FOL the derivation of the necessary lemmas could be done within the system itself. Of course, we don't get this for free in return we have to guide these systems very heavily in order to find the necessary proofs. Supplying lemmas to the Stanford Verifier lets us easily resort to well known results in literature, but we should always keep in mind that this practice bears the possible danger of introducing errors in the proof by supplying erroneous lemmas.

<sup>4.)</sup> building on [Fe77]

### 2.3. The Stanford verifier

This is a brief overview over the Stanford verifier, subsequently referred to as "the verifier". For more details the reader should consult [SV79]. The verifier accepts standard Pascal with minor syntactic changes. In addition certain specifications are required, some are optional. The verifier uses the (formally) documented program to generate a set of first order formulas (verification conditions) based on a set of Hoare style proof rules [HW73]. If all verification conditions are true, the program is consistent with its specifications. Note, that we distinguish "proving verification conditions" from "proving theorems and lemmas". The former is done completely mechanically by the verifier's prover while the latter refers to manual proofs. Occasionally we talk about "meta theorems". A meta theorem is a theorem about Hoare's logic itself. Examples are weak V-introduction and special proof rules about certain pointwise operations to be introduced in later chapters.

In the language accepted by the verifier the following documentation exists:

- entry and exit assertions are required for each procedure and the main program describing input and output predicates.
- An invariant is required for each loop.
- All loops constructed with jumps have to be cut by at least one assert statement, however, additional assert statements may appear at arbitrary places in the program.
- A comment statement is similar to an assert statement except that it does not break the path through the program where it is attached, rather it requires additional specification to be proven (see [SV79]).

The assertion language is quantifier free first order logic with arbitrary function and predicate symbols. All free variables in assertions have to be variables of the program. Special terms are added to express operations on complex data structures.

- $< A, [i], e >$  where  $A$  is an array denotes  $A$  with  $A[i]$  replaced by  $e$ .
- $< R, f, e >$  where  $R$  is a record with field  $f$  denotes  $R$  with  $R.f$  replaced by  $e$ .

A similar notation is used for operations on pointers [LS76]; this is discussed in chapter II in greater length.  
An important concept is the use of "virtual" code. The use of "auxiliary" variables in program verification has been pointed out in earlier work [LS75].

The concept of virtual code is much more general however. A virtual variable is one whose value at no point influences the value of "real" variables. We put no restrictions on the use of virtual variables. They may appear at any point where real variables appear, subject only to the restriction that they don't influence the outcome of a real computation. Furthermore we allow for virtual procedures, functions, types etc. It is shown in chapter IV how virtual code is useful in expressing the correctness of an algorithm in a natural way.

The verifier's built-in theorem prover contains complete decision procedures for various theories and their combination [NO78]. In addition it allows the user to specify "rules" axiomatizing free function and predicate symbols. Each rule expresses some logical axiom. Through different syntactic formats of rules heuristic information is conveyed to the prover how to use this particular fact in a proof. The logical language used for rules is essentially the same as for assertions within the program except that free variables are considered universally quantified.

The set of axioms and lemmas expressed by the rules required for a program proof is called the *logical basis* of this proof. Proofs are relative to their logical basis.

There are three different rule formats accepted by the verifier's theorem prover. Each rule expresses a certain logical fact. In addition different rule formats imply certain heuristics to the theorem prover. We briefly describe, the logical contents of the different rule formats; we are not concerned with the heuristic aspect of rules.

A replace rule<sup>5</sup> of the form

$\text{replace } T_1 \text{ where } F \text{ by } T_2$

is equivalent with  $F \Rightarrow T_1 = T_2$ . The where clause may be missing in which case  $F \equiv \text{true}$  is assumed. A rule of the form

$\text{infer } F_1 \text{ from } F_2$

from  $F_1 \text{ infer } F_2$

means  $F_2 \Rightarrow F_1$ . Similarly, the rule

$\text{from } F_1 \text{ infer } F_2$

means  $F_1 \Rightarrow F_2$ . The last two rule formats may appear with a whenever clause. Whenever clauses contain heuristic information for the theorem prover; they do not change the meaning of a rule.

5.) The name replace rule is misleading since the theorem prover does not implement a rewrite system.

### 3. Formal definition of programming languages

In order to be able to specify a compiler we need a method to formally define programming languages. Although there exists extensive literature on how to formalize language definitions this work seems not generally accepted; in most language designs informal English definitions prevail. The only exceptions are the by now generally accepted context free grammars. Many of the errors in current software can be attributed to a lack of specification. In the case of programming languages insufficient definitions leave ample room for major confusion of the user as well as the implementor. One result is the "try it out" attitude of many programmers who use the compiler rather than the manual as language definition. Also much of what is perceived as a compiler error is simply due to insufficient or ambiguous specification of languages.

In this section we briefly review various existing formal methods for the specification of programming languages. This section is far from being complete; rather the main emphasis is to motivate the particular choice of definitional formalisms used in the remainder of this thesis.

We adopt the common division of a language into syntax and semantics. The syntax describes how a program is externally represented as a string of characters while the semantics defines what a program means. Syntax is further subdivided into "micro syntax" or "token syntax" and the phrase structure of programs. Semantics is divided into "static" and "dynamic" semantics. Roughly, the former describes when a program is semantically valid; that is, it formalizes all type and declaration constraints. The dynamic semantics describes what the meaning of a semantically valid program is. Exactly what we mean by the meaning of a program depends very much on the particular formalism; we elaborate on this later. In addition to the components of a language definition outlined so far we introduce "tree transformations" which describe the relation between the phrase structure of a program and an "abstract" program which is input to the semantic analysis phase. This step may or may not be explicit in a particular formalism.

#### 3.1. Syntax

##### 3.1.1. Micro syntax

The external representation of languages is frequently described in terms of identifiers, numbers, reserved words and so on. The "micro syntax" of a language defines the syntax of these smallest syntactic entities, subsequently

called tokens.

A class of tokens (like identifiers) can in general be described as a regular language and in most textbooks a reference to this fact seems to settle this matter. But for a formal verification we have to be more specific and devise a precise definitional formalism.

Gries [Gr71] defines a scanner generator. The input to such a program can in some sense be considered to be the formal definition of the micro syntax. However, a more abstract theory is desirable and is developed in chapter III. To define the scanner formally we define (a) how the input string is to be broken into substrings which represent tokens and (b) how to map each of these strings into a token. Furthermore we define how these two mappings are connected to define the micro syntax.

In the course of this research several ways of giving a comprehensive definition of a scanner have been investigated. It turns out that there is one main problem in most methods: the definition tends to be very large and is very likely to be less comprehensible than a program implementing the scanner. This of course makes the usefulness of the whole verification questionable. For example, the definition of the micro syntax in terms of a finite automaton by specifying a transition table is extremely incomprehensible and error prone.

We adopt the following definitional schema to describe the micro syntax. We define a finite number of regular languages  $L_i$ , each of which characterizes a set of tokens. For example the set of identifiers, the set of numbers, the set of delimiters etc. With each language we associate a semantic function  $S_i$  which maps elements of  $L_i$  into the corresponding token. The name "semantic function" may be confusing in this context; it is to be understood in the sense that the meaning of a character string is a token. For example, a string of digits denotes a numeral.  $S_i$  are defined recursively in the style of a denotational definition. Finally, the scanner is given as a function  $scan$  which is defined in terms of  $L_i$  and  $S_i$ . An intuitive operational description of  $scan$  is

- discard all characters which are not an initial substring of any  $L_i$
- find the longest initial substring  $t$  of the input such that  $t \in L_i$  for some  $i$
- output the token  $S_i(t)$  and repeat this process until the input is exhausted.

##### 3.1.2. Phrase structure

Syntax is the best understood part of a language definition. In almost all cases the definition uses a context free grammar [AU72]. Various simplifying notations have been introduced, none of which alters the definitional power of

this formalism. The main deficiency of a grammar is that it merely defines the set of syntactically valid programs: it does not specify any particular output other than the derivation tree.

Several very efficient parsing methods are available. In particular for this compiler we use *LR* — *Parsing* [Kn65]. The key idea is to have a parsing table and a program interpreting this table. In some sense this interpreter defines a semantics for parsing tables and the tables themselves could be considered the formal definition of the syntax.

In chapter II we define the notion of a labeled context free grammar. For any unambiguous labeled context free grammar we define a mapping from token sequences to parse trees; labels of productions become labels in these trees. In chapter III the syntax of *LS* is defined in terms of these concepts.

### 3.1.3. Tree Transformations.

To provide a mapping of derivation trees into more suitable format we introduce tree transformations, functions mapping trees into trees. With this mapping the desired output format for the parser can be specified conveniently.

Tree transformations have been studied extensively in the literature [Kr74, Kr75, Ro71]. In particular DeRemer [De73, De74] studied their application in compiler construction. He views the whole translation process as a sequence of tree transformations. The starting point is a forest of singleton trees, the input characters and the result is the machine code. However, since transformations are purely syntactic they cannot serve as a formal definition of the semantics of the source language. A semantics is only given very indirectly by relating a program in source language to an (by definition) equivalent program in target language.

We define tree transformations using recursive functions on trees. We do not put any restrictions on these functions. It turns out, however, that for our source language the functions are of particularly simple nature and allow efficient implementation.

### 3.2. Semantics

Here the situation is completely different from that in the syntactic realm. Various different methods have been proposed, none of which appears to be generally accepted as a standard definition of semantics. We briefly outline several known methods for the definition of programming language semantics and motivate our particular choice of denotational semantics as a basis for compiler verification.

#### 3.2.1. Operational semantics

The most straightforward way to specify what a program does is to define an interpreter executing the program. Methods following this line are classified as operational definitions [Mc62, Mc63, Re72]. A special case of operational semantics is the Vienna definition method [LW69, We72].

A principal problem with operational semantics is that it does not define the meaning of programs abstractly. Rather, the "meaning of a program can only be determined for a particular set of input data by executing the program with the interpreter. It is not possible to reason about termination of programs on the basis of an operational semantics: if a program loops on a given set of input data then no will the interpreter. But maybe the most serious disadvantage is that an interpreter precludes many implementations of the language. Implementations that do not execute the program in the same way the interpreter does are very hard to verify since it amounts to proving the strong equivalence of two programs, the compiler and the interpreter.

On the other hand operational semantics has several important advantages. It provides a "complete" definition and it is rather easy to comprehend.

#### 3.2.2. Denotational semantics

In contrast mathematical or denotational semantics [Sc70, Sc71, Te76, Go79] introduces the meaning of a program as an object, a mapping from input data to results. Thus it is more abstract and allows reasoning about programs in a much more general way, e.g. programs can be compared for equivalence.

We argue that denotational semantics is best suited as a basis for a compiler proof. Firstly, it is the most abstract way of defining a language. As an important consequence a denotational definition relates very easily with other definition methods; other methods are subsumed and equivalence can easily be shown. Secondly, Scott provided a theory for the sound description of higher order functions and self application. Finally denotational definitions have been successfully applied in the definition of several realistic programming languages, among these are Pascal [Te77a], Algol 60 [Mor74], Algol 68 [Mi72c], SAL [MS76].

The semantics of our source language is defined analogously to that of Pascal in [Te77a]. In particular we use Tennent's separation of the semantics in static and dynamic semantics. However, there are some substantial differences in the way memory management is defined. Our definition is much more suitable as basis for a compiler proof.

### 3.2.3. Floyd - Hoare logic

In 1969 Hoare introduced an axiomatic system [Ho69] to reason about properties of programs. Since then various extensions to the method have been defined and theoretical results have been established. Of main interest are questions of soundness, completeness and the relation of an axiomatic definition to other methods [HL74, Do76, Co77, Ci79, GM79]. All these questions can only be answered with respect to a complementary definition of the language.

Although axiomatics is very well suited to reason about programs, only in exceptional cases can it serve as a formal definition of a language [Cl79, MG79]. In fact, the method gains much of its strength and usefulness because many subtle issues such as typechecking are purposely left out of consideration.

Even though it is conceivable to construct an axiomatic system which gives a complete detailed definition of all aspects of a language this seems not at all desirable. Attempts to extend axiomatics to completely define complex languages (of the order of Algol68) results in unnatural definitions [Sc78]; a semantics is only provided in a very indirect way.

It is unclear, how a compiler proof based on an axiomatic language definition can be given. A full second order logic seems necessary to prove statements of the form "all input output predicates which are satisfied by the source program are also satisfied by the target program".

### 3.2.4. Algebraic semantics

Algebraic semantics is based on the observation that the abstract syntax of a programming language is an initial algebra in a class of algebras  $C$ . Gouguen et al. [GT77] define any other algebra in  $C$  to be a possible semantics of the language. The relation between (abstract) syntax and semantics is given by the unique homomorphism between the initial algebra of  $C$  and any other member of  $C$ . The importance of the algebraic approach is that it unifies most other approaches to semantics. The main problem, however, is the construction of suitable semantics algebras. Denotational semantics is one way of specifying such an algebra.

Most of our presentation could be phrased in algebraic terms. However, we consider this unnecessary for our purpose and prefer not to do so.

An algebraic approach is also useful in the verification of compilers. In this case Cohn's central theorem [Co65, BL70] states that in a suitable algebraic framework (speaking in terms of syntax trees) it is only to prove that terminal nodes are translated correctly to ensure the correctness of the complete trans-

lator. This idea has been previously used to prove compiling algorithms (see section on previous work).

For our compiler proof the algebraic approach is not well suited. It requires to talk about the concept of homomorphism which is extremely hard, if not impossible to formalize suitably for present automated verification systems. Instead we use conventional computation induction which in some sense is equivalent (though it is a less abstract notion).

Further details go beyond the scope of this thesis, the reader should consult [Co65, Go78, GT77, MW72, Mo72, Mo73, Sc76a, Tw79, WM72].

### 3.2.5. Others

Some other methods for formalizing the semantics of programming languages have been proposed which cannot readily be classified in one of the above categories. Two relatively important ones are two-level or van Wijngaarden grammars [W169, W76] and attribute grammars [Kn68] both concepts are briefly described below.

The basic idea of van Wijngaarden grammars was to extend context free grammars to describe all context dependencies. However aspects of the dynamic semantics are still described informally in English. The main problems in using two-level grammars is that they are fairly difficult to relate to other formalisms. E.g. the definition of mode equality in [W76], section 7.3.1 and an equivalent first order definition of mode equality are very hard to relate and an equivalence proof is extremely tedious. A further problem is, that for two-level grammars there is no straightforward way of constructing a parser. Most existing Algol68 implementations use a standard context free grammar for their parser and ad hoc procedures to enforce context dependencies.

Another formalism is attribute grammars, first introduced by Knuth in [Kn68]. They have the distinct advantage that it is at least in principle possible to automatically generate an efficient compiler from a language definition in terms of attribute grammars. Several researchers have investigated this possibility and compiler generators are being constructed [Ka6a, Ka6b, Bo76, Cl79, De78]. The prime reason for not using this concept in this thesis is that it constitutes a description on a rather low level of abstraction. Also, an attribute grammar does not define a language; rather it relates it to another language (the machine code) by directly describing the translation process.

The connection between a set of attributes and an intuitive understanding of a language is very remote. In recent research Björner [Bj77] attempts to automatically generate an attribute grammar from a denotational definition.

This approach may in the future lead to efficient automatically constructed compilers; at this time, however, it is not useful for the research presented here.

### 3.3. Machine languages

In principle the techniques used for the source language definition are applicable to the target language as well. However, for the purpose of writing a compiler we are not so very much interested in "parsing" a given program in machine language. Rather, we perform the opposite step. But this is trivial and we omit the definition of a particular external representation for the target language. Consequently, all we are interested in is the semantics of the target language. As the reader may expect, we use denotational semantics for this purpose.

### 3.4. Summary

Altogether we define the source language by specifying a function *smean* mapping character strings into their meaning. The meaning of a program is simply a mapping from (input) files to (output) files. If our language had a more complicated file structure a different notion of meaning would be required.

The function *smean* in turn is defined as composition of functions defining the micro syntax, syntax, tree transformations, static, and dynamic semantics respectively.

The micro syntax is defined in terms of a function *scan*, mapping a string of characters into a string of tokens. The syntax is given by the function *parse* mapping token strings into a syntax tree. The function *tree<sub>r</sub>* maps syntax trees into abstract syntax trees. The static semantics is a function *ssem* defined on abstract syntax trees. It is the identity function if the program in question is semantically valid, otherwise *ssem* returns *error*. Finally, *dsem* maps an abstract syntax tree into its meaning, i.e. a function from files to files. Thus the function *smean* is given by

$$smean = dsem \circ ssem \circ tree_r \circ parse \circ scan.$$

The target language is defined by a function *tmean* mapping abstract code sequences into their meaning, again a mapping from input to output files.

The above functions are defined precisely in chapter III of this thesis. For now it suffices to assume their existence. We ignore the possibility of runtime and compile time errors for now.

### 4. Developing a verified compiler

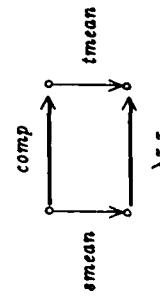
Informally the correctness of our compiler can be specified as follows:  
 "Whenever the compiler terminates without printing an error message it has generated correct code"

In this section we show how this statement can be formalised, given a formal definition of source and target language.

The above statement is all we prove about the compiler; we do not prove that it ever terminates, that it prints reasonable error messages, etc. Still, the above property is extremely important; it guarantees that the user of the compiler is never deceived into believing his program is translated correctly if it is not.

#### 4.1. What we prove

Source and target language are defined by the functions *smean* and *tmean* mapping a character string and a code sequence into their meaning. The compiler can be looked at as a function *comp* mapping a source program (a string of characters) into a target program (a code sequence). The compiler produces correct code if the following diagram commutes:



Where  $\lambda x.x$  is the identity function. In other words, the compiler is correct, if the meaning of the source program is equal to the meaning of the produced code. Although the above diagram is helpful in visualizing the correctness of a compiler it is not quite accurate in our case. The compiler to be written need not terminate for all input programs. Also, the compiler may issue meaningless error messages even if the input program is a valid LS program. If we consider printing of an error message as non-termination then the correct statement is given by that our compiler terminates then the above diagram commutes\*.

As described earlier, *smean* is given as a composition

$$smean = dsem \circ ssem \circ tree_r \circ parse \circ scan.$$

According to this definition the compiler is divided into four (*treeir* and *parse* are computed in one program module) separate programs executed sequentially. Each program has an input domain and an output domain, such that output and input of two subsequent programs match.

The individual program modules can be thought of as procedures with the following specifications:

```
procedure sc(in: char sequence; var out: token sequence);
```

```
exit out = scan(in)
```

```
procedure ps(in: token sequence; var out: abstract syntaxtree);  
exit out = treeir(parse(in))
```

```
procedure ss(in: abstract syntaxtree; var out: abstract syntaxtree);
```

```
exit out = ssem(in)
```

```
procedure cg(in: abstract syntaxtree; var out: code sequence);  
exit out = dem(in)
```

The compiler consists of the sequential execution of these four programs.

This can be written as

```
sc(in, tks); ps(tks, as1); ss(as1, as2); cg(as2, out).
```

After executing this sequence it is true that

```
tsem(out) = dsem(ssem(treeir(parse(scan(in)))))
```

By the standard interpretation for Hoare formulas this means that if the compiler terminates condition (\*) holds. But (\*) is just another way of stating that the diagram given above commutes.

A precise specification of the compiler and its components will be given in chapter III.

In the next subsection some principal techniques used in the verification of the compiler are discussed. Following this, we describe the concepts used in the implementation and proof of the individual components of the compiler.

#### 4.2. Representation

Objects manipulated by a computer program are bits. If we are talking about integers, booleans, characters, or even sets or queues we actually are talking about a representation of these objects in terms of bits. For example, an integer number may be represented by a sequence of 32 bits; the value

of such a sequence is derived by interpreting the sequence as a binary representation. Alternatively, if we want to talk about arbitrarily large integers, we may have a linked list of elements of a base 1000 representation.

In general an abstract object is represented by specifying a data structure that is used to store such an object and by giving a mapping from this data structure to the desired object. To represent sets for example, we may choose linked list as representation and define a "representation function" as

$$\text{setrep}(l) = \text{if } \text{null}(l) \text{ then } \{\} \text{ else } \{\text{car}(l)\} \cup \text{setrep}(\text{cdr}(l)).$$

In chapter II we expand on this idea. In particular we are interested how function domains and functions on functions can be represented. In many cases the representation theory can be used to automatically select suitable data structures. A possible application are compiler generators, where data structures could be derived from domain definitions of the denotational semantics.

The idea of representation is relevant for all implementations [J079] and proofs. It significantly simplifies the proof of a refinement step, given a proof for the original program. The following trivial example explains this situation. Suppose we are given an abstract procedure *add* with the entry *exit specification*

$$\{A = A_0\} \text{add}(A, b) \{A = A_0 \cup \{b\}\}.$$

Given the simple theorem  $(A \cup \{b\}) \cup \{c\} = A \cup \{b, c\}$  we can prove the correctness of the program

$$\{A = A_0\} \text{add}(A, b); \text{add}(A, c) \{A = A_0 \cup \{b, c\}\}.$$

Suppose now, that we implement sets as linked lists as above, then *add* operates on a list *l* instead of a set *A*. However, the documentation is still written in terms of sets. I.e. the refined version of *add* can be specified as

$$\{\text{setrep}(l) = A_0\} \text{add}(l, b) \{\text{setrep}(l) = A_0 \cup \{b\}\}.$$

Using the refined version of *add* does not change the structure of the proof of the using program. The only change is, that the documentation in the using program has to be changed to account for representation functions. The above example becomes

$$\{\text{setrep}(l) = A_0\} \text{add}(l, b); \text{add}(l, c) \{\text{setrep}(l) = A_0 \cup \{b, c\}\}.$$

No additional theorems are necessary to prove this program. Only in proving the implementation of *add* do we need properties of *setrep*.

#### 4.3. Scanner

For our particular scanner we prove a theorem stating that for any two languages  $L_i$  and  $L_j$ , their respective sets of initial substrings are disjoint. This is useful for an efficient implementation. For example seeing a letter in the input string means that this must be the beginning of an element of the set of identifiers and cannot belong to any other  $L_i$ .

In addition we need several trivial theorems about sequences of characters and tokens. They are of the nature of

$$\begin{aligned} \langle \rangle z &= z \cdot \langle \rangle = z \\ xy &= xz \rightarrow y = z \end{aligned}$$

and deserve no further attention.

The implementation follows the recursive definition of  $\text{scm}$  and  $S_i$  and uses standard recursion removal techniques. In addition the above theorem about  $L_i$  allows several shortcuts to improve efficiency.

#### 4.4. Parser

We use a LR parser [AJ74] to construct the required parse tree. We are not concerned with LR table construction; rather we assume that we have a correct table generator at our disposal.

The theory of LR parsing is fairly well understood, but a formal machine checked proof of an LR parser has not yet been published. The main problem we were faced with was to precisely define the properties of parsing tables required to prove the parser correct. More details are presented in chapter IV.

Again, of course, the parser heavily uses the representation theory and sequences as well as trees.

We talked so much about a stepwise development of program so the reader may be interested in some statistics. We have written and proven 10 successive versions of the parser. The starting point was a barebone version which assumed the correctness of parsing actions and did not produce output. The implementation of parsing actions followed. Next, in several versions we introduced the construction of parse trees and finally added the tree transformation step to produce abstract syntax trees.

The parser is a good example to illustrate the use of virtual code. It would of course be possible to have the parser construct a parse tree and after the program is parsed to apply tree transformations to produce the ultimate output. Instead, the parser produces a parse tree and an abstract syntax tree in parallel; part of the proof is to show that the abstract syntax is the image of

the parse tree under the tree transformations at any intermediate stage. We then can prove that the parser satisfies the following specifications

$$\{\text{input} = \text{input}_0\} \text{pa}(\text{pt} = \text{parSet}(\text{input}_0)) \wedge \text{ast} = \text{tree}(\text{pt})\}.$$

It turns out, that the abstract syntax tree ( $\text{ast}$ ) at no point depends on the parse tree ( $\text{pt}$ ); rather it is derived directly from the parsing actions performed. Thus, the parse tree can be virtual code; it never has to be actually computed. The exit specification of the parser then implies that  $\text{ast} = \text{parse} \circ \text{tree}(\text{input}_0)$ .

#### 4.5. Semantic analysis

The static semantics is defined by recursive functions in Scott's logic of computable functions. In principle the implementation is straightforward but there are two theoretical problems.

The first problem is to express the definition in the assertion language accepted by the verifier. This is complicated by the fact that we have to deal with higher order functions (i.e. function mapping functions into functions) and that we have to find a consistent treatment of undefined elements ( $\perp$ ).

The second problem is what to do with the least fixed point operator ( $\text{fix}$ ) that is used in the definition. If  $\text{fix}$  is used to define a recursive function then the solution is, of course, to compute this function recursively. The interesting case is, however, when  $\text{fix}$  is used to describe a data object which has to be computed.

The first problem can be solved fairly easily using the concept of representation functions; details are presented in chapter II.

The key idea to solve the second problem is referred to as "operationalisation" of fixed points. Let  $f$  be a function from  $D$  into  $D$ , we write  $f \in D \rightarrow D$ . The least fixed point of this function  $\text{fix } f$  is an element of  $D$ .<sup>6</sup> Given a representation  $r$  for elements in  $D$  and a representation of  $f \in D \rightarrow D$  the problem is to compute the representation of  $\text{fix } f$ . Note, if we are allowed to choose a suitable representation, this problem has always a solution. Just take the symbolic representation  $f$  to represent the function  $f$ , then  $\text{fix } f$  is a symbolic representation of the fixed point. Of course, we are interested in a solution to this problem for more efficient representations.

Let us look at a simple example. Suppose the domain  $D$  is the set of all finite and infinite lists. As a representation we choose Lisp lists and define

$$\text{rep}(l) = \text{if null}(l) \text{ then } () \text{ else } (\text{car}(l)) \text{ rep}(\text{cdr}(l)).$$

<sup>6</sup>) Precise definition of these notations can be found in chapter II.

Under  $\text{rep}$  a circular Lisp list represents an infinite abstract list. Clearly, not all infinite lists can be represented in this way, but a representation is not required to be surjective.

Let us now consider how we can compute the fixed point  $\text{fix}(\lambda z.(a \cdot z))$  which is a list consisting of infinitely many  $a$ 's. First, we need a representation for functions from lists to lists. We take a pair of two pointers (S-expressions) for this purpose and define

$$\text{funcrep}(p, q) = \lambda t. \begin{cases} \text{If } \text{null}(p) \text{ then } () \\ \quad \text{If } \text{eq}(p, q) \text{ then } t \text{ else } (\text{car}(p)) \cdot \text{funcrep}(\text{cdr}(p), q)(t) \end{cases}$$

To represent  $\lambda z.(a \cdot z)$  we simply take an arbitrary S-expression  $s$  other than  $\text{nil}$  and use the pair  $(\text{cons}(a, s), s)$ . Given an arbitrary pair  $(p, q)$ , one can compute  $p$  such that  $\text{rep}(p) = \text{fix}(\text{funcrep}(p, q))$  as follows. Find a sublist  $t$  of  $p$  such that  $\text{eq}(\text{car}(t), q)$ ; compute  $\text{rep}(\text{cdr}(t, p))$ ; let  $\hat{p} = p$ . It is easy to prove that after these steps  $\hat{p}$  has the desired property.

Applied to our example, we let  $p = \text{cons}(a, s)$  and construct  $\text{rep}(p, p)$ . Clearly,  $p$  now points to a cyclic list which under  $\text{rep}$  represents an infinite list of  $a$ 's.

In chapter II we formalize the above argument and prove a general theorem that then enables us to efficiently compute the least fixed points required for the semantic analysis.

Given the above theory, the implementation of the semantic analysis phase is straightforward.

#### 4.6. Code generation

##### 4.6.1. Necessary theory

We choose a fairly low level definition of the source language. For example, addressing is done relative to base addresses. For new procedure invocations new base addresses are set up. A display mechanism describes how memory of outer scopes is to be accessed. These concepts are familiar to the compiler writer and are well suited as our specification.

For some statements of  $L_S$  we give a modified definition which is better suited for the verification of the code generation. For example, the definition of the while loop is originally given by a least fixed point. In the revised definition the meaning of the while loop is described by a conditional and explicit jumps. We prove manually that the revised definition is equivalent to the original one.

Also, we give a revised definition of the target language. We add the concept of blocks to the target language. All a block does is to hide labels defined inside from the outside. This concept is advantageous in generating code for example for the while loop. The necessary code requires introduction of labels in the target language that are not existent in the source language; it is convenient to be able to hide these in a block.

But we do not generate code for a different language. Rather, we prove that if all labels in a target program with blocks are distinct, then we can omit all blocks without changing the meaning of the program. The code generation produces "virtual" code containing blocks and real code that is identical except that blocks are omitted. The specification of the code generation is given as

$$\{ \text{ast} = \text{ast}_0 \} \text{ cg}$$

$\{ t. \text{mean}(\text{code}_1) = \text{dsem}(\text{ast}_0) \wedge \text{code}_2 = \text{flat}(\text{code}_1) \wedge \text{distinct}(\text{code}_1) \}$

$\text{code}_1$  contains blocks and is proven correct.  $\text{code}_2$  is generated without blocks but is otherwise identical to  $\text{code}_1$ , expressed as  $\text{code}_2 = \text{flat}(\text{code}_1)$ . Finally, we prove that all labels in  $\text{code}_1$  are distinct. Thus the above theorem is applicable and we may conclude  $\text{tmean}(\text{code}_2) = \text{dsem}(\text{ast}_0)$ . Observe, that this reasoning is analogous to that in the parser.

##### 4.6.2. Implementation

With the above theory at hand the actual implementation is fairly simple. We have a set of procedures to generate fixed code sequences for primitive operations such as updating a location, accessing the value of a location and so on.

Code for more compound entities is generated recursively. Let us illustrate this with a very simple example. Suppose the recursive procedure  $\text{code}$  generates code for expressions. We have the specification

$$\{ \text{true} \} \text{code}(E, z) \{ "z \text{ computes } E, \text{ puts result on stack}" \}$$

Inside  $\text{code}$  we have a branch which deals with binary operations:

```
if  $E = "E_1 + E_2"$  then
begin
  code( $E_1, z_1$ );
  code( $E_2, z_2$ );
  primcode(+, z_3);
   $z \leftarrow z_1 \cdot z_2 \cdot z_3$ ;
end
```

`primicode` is a procedure generating code for primitive operations; in this case for “+” it produces code that takes the two top elements of the stack and pushes back their sum. Although this example is very much simplified, it demonstrates the basic principle of the code generation and its proof.

The reader may wonder about one more problem: how are fixed points of the semantic definition handled in the code generation? The solution is simple in this case. The meaning of a loop in the source language is described as a fixed point, say  $fiz\ s$ . For this loop we generate code whose meaning in the target language is defined as a fixed point as well, say  $fiz\ t$ . To show that both fixed points are equal, we merely have to prove that  $t = s$ .

## 5. Related work

### 5.1. Previous work on compiler verification

Correct compilers have been of great interest to researchers from the very beginning of verification ideas. We summarise the work that has been done in this area and give a short characterization of the particular innovations introduced in each case.

The first attempt to verify a compiler or rather a translation algorithm, has been undertaken by McCarthy and Painter [MP66] in 1966 (see also [Pa67]). Though limited in scope this work has been of great impact and many researchers have directly built upon it. The language considered are arithmetic expressions. The compiler is a set of recursive functions and the proof proceeds by recursion induction.<sup>7</sup> The underlying semantical definition is operational.

A machine checked proof of the McCarthy/Painter compiler has been done by W. Diffie but has not been published.

Kaplan [Ka67] has built on McCarthy/Painter and proved a compiler for a yet more complex language including assignments and loops, still using the same proof technique. Again in [Bu69] R. Burstall proves the McCarthy/Painter compiler by structural induction. But as he points out, the proof is very similar to one given by recursion induction.

A completely new aspect enters the field of compiler verification in 1969 when Burstall and Landin [BL69] apply an algebraic method to greatly reduce the amount of work required in proving a compiler.

<sup>7.)</sup> different recursion principles are discussed in chapter II.

The same proof method is used by Milner and Weyhrauch [MW72, WM72] in 1972 to verify a compiler for an algolic language. But they introduce yet another innovation: semantics and proof are described in Scott's logic of computable functions. F.L. Morris uses the same algebraic method in [Mo72, Mo73]. More recently Thatcher, Wagner, and Wright [TW79] put forth similar ideas.

The proof techniques used by McCarthy and Painter are again employed by London [Lo71, Lo72] to prove an existing compiler for a simple version of LISP.

In recent years some authors tried to use axiomatic semantics as a basis for a compiler proof [Ch76, CM75]. But there appear to be some difficult problems. There is no notion of semantics in these proofs at all. Source and target languages are “defined” by a set of rules. In addition, static semantic issues are totally ignored. Again, the language handled is not much more interesting than that of the McCarthy/Painter compiler.

Lynn [Ly78] is the first to consider the problem of user defined functions in a compiler for LISP. He too uses Hoare style proof rules to specify the semantics of source and target language.

Yet a different approach is taken by Boyer and Moore. Their LJSP theorem prover [BM77] has been shown to be capable of verifying an optimizing compiler for LISP expressions. Their system is capable of automatically synthesizing induction hypotheses for suitable recursive functions.

Two researchers continued to use denotational semantics as a basis for a compiler proof. First there are Milne and Strachey which in their book [MS76] give a proof of a compiling algorithm for a language of the complexity of Aigo68. The proof is given completely by hand and the target language is a hypothetical machine language. Avra Cohn [Co79a, Co79b] proves the correctness of several components of a compiling algorithm using the Edinburgh implementation of LCF [GM75]. One problem she is interested in is the compilation of recursive procedures into stack implementations. Her main emphasis, though, is on the technique of using LCF and automating proofs in LCF.

A totally different approach is taken by H. Samet [Sa75]. Instead of proving the correctness of a compiler he proves equivalence of a particular source program and the produced code. Although this approach is in general much more expensive than verifying a compiler, it is useful in situations where a correct translation is imperative and a correct compiler is not available.

### 5.2. Relation to our work

Most practitioners would not consider "compilers" verified previously to be "real compilers". Our emphasis is to verify a compiler that is - at least in principle - a practical, usable compiler. Therefore we choose a realistic source language in which a programmer might want to actually program. The only language comparable in size and complexity is *SAL*, considered by Milne and Strachey in [MS76]. But they consider a very abstract compiler only and their attention is limited to code generation, all proves are manual.

Milne and Strachey use a semantic definition of a higher level of abstraction than the definition of *LS* used in this work. If we were to use a more abstract definition a lower level definition would have to be developed as part of the theory for the code generation. The techniques and results in [MS76] are directly applicable for the necessary equivalence proof. In this sense the work in [MS76] is complementary to ours.

Cohn's work represents an important step towards mechanically checking proofs of the kind we give manually to derive theorems necessary for the code generation.

For our compiler to be practical we consider not only the code generation part but also verify a scanner, parser, and static semantic analysis. These issues have not been dealt with previously.

Although the development of a suitable logical basis for our compiler verification is done manually all actual program proofs are completely machine checked. No comparable mechanical verification attempt is known to the author.

An important part of this research deals with technical issues that are prerequisites for a mechanical compiler verification. These problems have not been considered before in this form. Some main points are

- formalization of compiler correctness suitable as input for a mechanical verification system,
- structuring of programs such that mechanical verification becomes manageable, and
- development of specialized verification techniques to deal with pointer operations, least fixed points, and quantification.

### 5.3. Compiler generators

It would, of course, be much more efficient could we verify a compiler generator. We then could generate correct compilers for any specification.

But, given the current state of the art, this goal is somewhat unrealistic. We do not even know how to write compiler generators rather than how to verify them. Though much has been written about this subject no practical system exists to my knowledge. Most system called compiler generators turn out to be just parser generators.

Important work towards automatic generation of compilers has been done though. We should mention P. Mosses' thesis [Mo73b] and the work based on attribute grammars [Ka76a, Ka76b, Bo76, CH79, Bj77].

Our work is not totally unrelated to compiler generation. Firstly, the parsing algorithm implemented operates table driven, thus it remains correct for any language that can be defined by an LR(1) grammar.

Further, we show how a program and its proof can be systematically derived from a denotational semantics. These results give some indication as to how this process or at least some parts of it can be automated.

## 6. Organization of this thesis

In chapter II we introduce notations and relevant theories required for the compiler proof. Among other things we consider

- first order logic and its application as assertion language
- Scott's logic of computable functions
- Axiomatization of Scott's theory in logic
- Proofs in logic and Scott's theory and their relation
- Theory of representation
- Treatment of fixed points

In chapter III we describe the source and target languages *LS* and *LT*. First, we give an informal introduction to the source language followed by the formal definition of scanner, parser, static and dynamic semantics. For each of these parts we define the definitional formalism used. We introduce a hypothetical machine executing the target language *LT*. A formal definition of *LT* is given.

Chapter IV describes precisely how the correctness of the compiler and its components is specified in a format accepted by the verifier. We then discuss the systematic development of the individual components and demonstrate the proof principles employed.

In the concluding chapter V we summarise our results and consider possible changes, improvements and extensions of the compiler. We discuss implications of this work on the design of verification systems and more sophisticated program development systems that encompass the whole design cycle of programs.

## Chapter II. Theoretical framework

### 1. Basic concepts

#### 1.1. Functions

Let  $A$  and  $B$  be two arbitrary sets;  $A \rightarrow B$  denotes the set of all total functions from  $A$  to  $B$ . " $\rightarrow$ " associates to the right. We assume that all functions are either constants or have exactly one argument. This allows to write function application in "curried" form as  $f x$ . This "juxtaposition" associates to the left; i.e.  $f x y$  means  $(f x)y$ . Functions with several arguments as  $f \in A \times B \rightarrow C$  are applied to tuples as in  $f(x,y)$ . Alternatively,  $f$  can be defined for the isomorphic domain  $A \rightarrow (B \rightarrow C)$  in which case we write  $f x y$ . We use the semicolon ( $:$ ) to break the normal precedence, i.e.  $f; g h$  means  $f(g h)$ . Several terms separated by ; associate to the right, i.e.  $T_1; T_2; T_3$  means  $T_1(T_2(T_3))$ .

New functions can be constructed from old ones by using composition, conditionals, and  $\lambda$ -abstraction as defined below.<sup>1</sup> This language is essentially Church's  $\lambda$ -calculus [Ch5]. However, we are not so much interested in the  $\lambda$ -calculus as a formal system; rather we use it as a notational vehicle.

[ $\lambda$ -abstraction] Let  $T(x)$  be a term with the potentially free variable  $x$ , then  $\lambda x.T(x)$  denotes a function  $f$  such that  $f y = T(y)$ .

The notations "let  $x = T_1$  in  $T_2$ " and " $T_2$  where  $x = T_1$ " are both equivalent to  $(\lambda x.T_2)T_1$ .

Some special sets used are integer  $N$  and truthvalues {TT, FF}.

1.) We are not strictly formal here, rather we identify terms denoting functions with the function denoted. A formal treatment would require to introduce the semantics of the  $\lambda$ -calculus (e.g. [B177]).

**[Conditionals]** The conditional  $\text{if } T_1 \text{ then } T_2$  else  $T_3$  has the usual meaning; it is only well defined if  $T_1$  denotes an element in  $\{TT, FF\}$ . This definition of conditionals will be extended to domains in the next section.

**[Redefinition]** Let  $f \in A \rightarrow B$ , then  $f[z/y]$  denotes the function

$$\lambda\epsilon. \text{if } \epsilon = y \text{ then } z \text{ else } f\epsilon.$$

### 1.2. First order logic

We assume that the reader is familiar with basic concepts of logic (see e.g. Mendelson [Me64] or [KlG7]). The following definition of first order logic is written in the style of denotational definitions and is mainly included here to give the reader an idea of the flavour of these definitions using a familiar example. Note, however, that all sets and functions are ordinary mathematical sets and total functions and do not have any added structure as they do in denotational semantics. This is possible since the domains in question are very simple and the problem of self application does not arise. Parentheses  $[x]$  are used to indicate that  $x$  is a syntactic object; this convention will be obeyed throughout this thesis.

#### 1.2.1. Syntax

$x \in X$   
 $f_i \in Y_i$   
 $p_i \in P_i$   
 $t \in T = x \mid f_n(t_1, \dots, t_n) \mid t_2 \mid p_n(t_1, \dots, t_n) \mid u_1 \wedge u_2 \mid \neg u \mid \forall z. u$

Terms  $u \in F = \text{true} \mid \text{false} \mid t_1 = t_2 \mid p_n(t_1, \dots, t_n) \mid u_1 \wedge u_2 \mid \neg u \mid \forall z. u$  Formulas  $t$  are abbreviated as  $f_i$ ; if no ambiguities arise the arity index is omitted. Parentheses may be used arbitrarily to disambiguate. The formulas  $u_1 \vee u_2$ ,  $u_1 \Rightarrow u_2$ , and  $\exists z. u$  are abbreviations for  $\neg(\neg u_1 \wedge \neg u_2)$ ,  $\neg u_1 \vee u_2$ , and  $\neg(\forall z. \neg u)$  respectively. Some common function symbols like  $+$ ,  $-$ , etc. are written infix.

#### 1.2.2. Semantics

The semantics is defined by specifying the universe  $U$  of interpretations and ‘meaning’ functions mapping terms into  $U$  and formulas into truth values. Formally we write:

$$U \quad \text{Universe}$$

$\text{Fun}_i = U \rightarrow U \rightarrow \dots \rightarrow U$   
 $\text{Pre}_i = U \rightarrow \dots \rightarrow U \rightarrow (TT, FF)$   
 $\phi \in \Phi = (X \rightarrow U) \times (Y_i \rightarrow \text{Fun}_i) \times (P_i \rightarrow \text{Pre}_i)$

The last line, for example, is to be read as: “an interpretation  $\phi$  is a function that maps variables into the universe  $U$ ,  $i$ -ary function symbols into  $i$ -ary functions and  $i$ -ary predicate symbols into  $i$ -ary predicates.”

$$\boxed{T \in T \rightarrow \Phi \rightarrow U}$$

$$\begin{aligned} T[z]\phi &= \phi[z] \\ T[t_1, \dots, t_n]\phi &= (\phi[t_1])(T[t_1]\phi) \dots (T[t_n]\phi) \end{aligned}$$

$$\boxed{\mathcal{I} \in F \rightarrow \Phi \rightarrow (TT, FF)}$$

$$\begin{aligned} \mathcal{I}[\text{true}]\phi &= TT \\ \mathcal{I}[\text{false}]\phi &= FF \\ \mathcal{I}[t_1 = t_2]\phi &= (T[t_1]\phi = T[t_2]\phi) \\ \mathcal{I}[p(t_1, \dots, t_n)]\phi &= (\phi[p])(T[t_1]\phi) \dots (T[t_n]\phi) \\ \mathcal{I}[u_1 \wedge u_2]\phi &= \text{if } \mathcal{I}[u_1]\phi \text{ then } \mathcal{I}[u_2]\phi \text{ else } FF \\ \mathcal{I}[f \neg u]\phi &= \text{if } \mathcal{I}[u]\phi \text{ then } FF \text{ else } TT \\ \mathcal{I}[\forall z. u]\phi &= (\lambda\epsilon. \mathcal{I}[u]\phi[e/z]) = \lambda\epsilon. \mathcal{I}[u]\phi \end{aligned}$$

Other logical notions such as validity, satisfiability, proof, etc. can be defined in the obvious way in our formalism.

#### 2. Scott's logic of computable functions

The mathematical theory underlying denotational semantics has been developed by Dana Scott [Sc70]. In lack of a better name we refer to this theory as Scott's logic of computable functions or Scott theory for short. The theory we use differs slightly from that originally proposed by Scott; we assume an interpretation over complete partial orders (cpo's) [P78] rather than over lattices.

In this section we briefly summarize the relevant concepts and definitions of Scott's logic of computable functions. For a more complete treatment of this subject see [S77].

### 2.1. Basic definitions

**[Quotients]** Let  $S$  be a set and  $R$  be an equivalence relation on  $S$ , then the quotient  $S/R$  is the set of equivalence classes of  $R$ . If  $f$  is a function defined on  $S$ , then  $S/f$  is  $S/R$  where  $R = \{(x, y) \mid f(x) = f(y)\}$ .

**[Directed set]** If  $(D, \sqsubseteq)$  is a partial order, then  $S \subseteq D$  is directed if  $S$  has an upper bound.

**[CPO]**  $(D, \sqsubseteq)$  is a complete partial order (CPO) if  $(D, \sqsubseteq)$  is a partial order with

- bottom ( $\perp_D \in D$ ) is the least element in  $D$ , and
- every directed subset  $S$  of  $D$  has a least upper bound ( $\sqcup_S$ ) in  $D$ .

**[Domains]** For the purpose of this thesis a domain is a CPO; we write  $D$  for a domain and omit the ordering  $\sqsubseteq$  if it is clear from the context.<sup>2</sup>

Some domains have a set of error elements  $E$ ; i.e. they are of the form  $D + E$ . We sometimes omit  $E$  in the definitions. Operations which are defined on  $D$  extend to  $D + E$  in a strict way, that is they are the identity on  $E$ . If it is not necessary to distinguish different errors we use  $E = \{\}\!$ .

**[Flat Domains]** Given an arbitrary set  $S$  the CPO  $(S_\perp, \sqsubseteq)$  with  $S_\perp = S \cup \{\perp\}$  and  $a \sqsubseteq b$  iff  $a = \perp$  is called a Flat Domain. Elements in  $S$  are called proper.

Flat domains are important because they allow us to construct a domain from a given set. Some typical examples we will use later are

$$T = \{TT, FF\}_\perp, \quad N = \mathbb{Z}_\perp$$

where  $\mathbb{Z}$  are the usual integers.

### 2.2. Operations on domains

**[Sum]** If  $D_1$  and  $D_2$  are two domains, then  $D_1 + D_2$  is the sum of  $D_1$  and  $D_2$ . Elements of  $D_1 + D_2$  are elements from  $\{(1, z) \mid z \in D_1\} \cup \{(2, x) \mid x \in D_2\}$

<sup>2.)</sup> Topological issues such as finite basis are irrelevant for our discussion. For more details see [Sci72a, Sci76b, Pfl78].

with  $(1, \perp_{D_1}) = (2, \perp_{D_2}) = \perp_{D_1 + D_2}$ .  $\sqsubseteq$  extends to the sum according to  $(i, x) \sqsubseteq (j, y)$  iff  $i = j \wedge x \sqsubseteq y$ .

If  $D = D_1 + D_2$  and  $x \in D$ , then  $x.D$  denotes  $(i, x)$ , the injection of  $x$  into  $D$ . For  $z = (i, y) \in D$  the predicate  $z \in D_j$  is true iff  $i = j$ . If  $z = (i, y) \in D$ , then  $z|_{D_j}$  is the projection of  $z$  onto the component  $D_j$ , i.e.

$$\langle i, y \rangle|_{D_j} = \begin{cases} \perp & \text{if } i \neq j \\ y & \text{if } i = j \end{cases}$$

Frequently the explicit projections and injections are omitted if no ambiguities arise.

**[Product]**  $D_1 \times D_2$  is the Product domain of  $D_1$  and  $D_2$  with  $(z_1, z_2) \sqsubseteq (y_1, y_2)$  iff  $z_1 \sqsubseteq y_1 \wedge z_2 \sqsubseteq y_2$ . If  $z \in D_1 \times D_2$  then  $z = (z^{#1}, z^{#2})$ . Note, while  $z_i$  is a variable, "supercript #i" is a projection function defined on product domains.

Sum and product spaces extend to more than two dimensions in the obvious way.

**[Continuous]** A function  $f \in D_1 \rightarrow D_2$  is continuous if for all directed subsets  $S \subseteq D_1$  we have  $f(\sqcup S) = \sqcup f(S)$ , where  $f(S) = \{f(s) \mid s \in S\}$ .

**[Function domains]**  $D_1 \rightarrow D_2$  is the domain of continuous functions from  $D_1$  into  $D_2$  ordered by  $f \sqsubseteq g$  iff  $\forall x \in D_1, f(x) \sqsubseteq g(x)$ .<sup>3</sup>

**[Monotonic]** A function  $f$  is monotonic if  $\mathbf{x} \sqsubseteq \mathbf{y} \rightarrow f(\mathbf{x}) \sqsubseteq f(\mathbf{y})$ .

**Corollary:** Every continuous function is monotonic. (Consider the directed set  $S = \{x, y\}$  with  $\mathbf{x} \sqsubseteq \mathbf{y}$ .)

**[strict]** A function  $f \in D_1 \rightarrow D_2$  is strict if  $f(\perp_{D_1}) = \perp_{D_2}$ .

**Corollary:** Any strict function defined on a flat domain is monotonic and continuous. (All directed sets of a flat domain are trivial.)

<sup>3.)</sup> Note the difference between  $A \rightarrow B$  and  $A \rightarrow B$ .

[Strict product] Let  $R$  be the relation on  $D_1 \times D_2$  such that

$$(z_1, z_2)R(y_1, y_2) \text{ iff } (z_1 = \perp_D, V z_2 = \perp_{D_2}) \wedge (y_1 = \perp_{D_1}, V y_2 = \perp_{D_2}).$$

The domain  $D_1 \times_0 D_2 = (D_1 \times D_2)/R$  is the Strict Product of  $D_1$  and  $D_2$ .

Arbitrary expressions over domains using the operators  $+$ ,  $\times$ ,  $\rightarrow$  define new domains. We use the convention that  $\rightarrow$  associates to the right.  $+$  and  $\times$  are associative (up to isomorphism).

Equations defining new domains may be recursive (mutually). Domains defined recursively are well defined as is shown in [Sc76b].<sup>4)</sup>

As an example consider lists  $L$  over a domain of atoms  $A$ . The domain  $L$  can be defined recursively as  $L = A + A \times L$ . Let us briefly investigate the structure of  $L$ . Clearly,  $L$  contains the chain

$$\perp \sqsubseteq \langle a, \perp \rangle \sqsubseteq \langle a, \langle a, \perp \rangle \rangle \sqsubseteq \langle a, \langle a, \langle a, \perp \rangle \rangle \rangle \sqsubseteq \dots$$

For  $L$  to be a cpo it has to contain the limit of this chain which is an infinite list. This is undesirable if, for example, we are only interested in finite lists as in LISP.

We get a different situation if we define  $L = A + A \times_0 L$ , using the strict product. Now we have

$$\perp = \langle a, \perp \rangle = \langle a, \langle a, \perp \rangle \rangle = \langle a, \langle a, \langle a, \perp \rangle \rangle \rangle = \dots,$$

thus at least this sequence does not generate infinite objects. In fact, it can be shown that  $L$  contains exactly the finite lists over  $A$  together with  $\perp$ ;  $L$  is the flat cpo of finite lists over  $A$ .

Another example for domains with infinite objects are binary trees  $B$  with  $B = B \times B + A$ . Again,  $B$  contains infinite objects. Using strict products gives us exactly the finite binary trees.

[Fixed points] Let  $f \in D \rightarrow D$ , then  $\text{fix } f$  denotes the least fixed point of  $f$ , that is the least element of  $\{g \mid g = f(g)\}$ .

For monotonic  $f$  the least fixed point exists; for continuous  $f$  the least fixed point is given by  $\text{fix } f = \sqcup\{f_i\}$  where  $f_0 = \perp$ ,  $f_{i+1} = f(f_i)$ .

### 2.3. Conditionals

Conditionals **if** — **then** — **else** appearing in formulas of Scott's logic

<sup>4.)</sup> The proof is based on a universal domain  $U = U \rightarrow U$ , other domains are defined as images of  $U$  under retracts [Sc76b].

are defined completely strict in the condition.

$$\text{if } E_0 \text{ then } E_1 \text{ else } E_2 = \begin{cases} \perp & \text{if } E_0 = \perp \\ E_1 & \text{if } E_0 = \text{true} \\ E_2 & \text{if } E_0 = \text{false} \end{cases}$$

The **if** construct may appear without an **else** part in which case it is defined as:

$$\text{if } E_0 \text{ then } E_1 = \begin{cases} \perp & \text{if } E_0 = \perp \\ E_1 & \text{if } E_0 = \text{true} \\ \perp & \text{if } E_0 = \text{false} \end{cases}$$

Both versions of the conditional are continuous.

### 2.4. Lists

Finite lists (or sequences) over a given domain occur so frequently that we introduce a special operator. If  $D$  is a domain then  $D^*$  is defined as  $D^* = D + D \times_0 D^*$ .

We write  $\langle \rangle$  for the empty sequence.  $\langle z_1, \dots, z_n \rangle$  is a sequence with elements  $z_1, \dots, z_n$ . The function  $\text{hd}$  and  $\text{tl}$  yield the head and tail of a sequence and  $\perp$  if applied to  $\langle \rangle$ .

A centered dot “.” is used to denote concatenation. When no ambiguities arise we write  $\langle x \rangle w$  as  $x \cdot w$ .

If  $s$  is a sequence  $|s|$  denotes the length of  $s$ .

Let  $f \in D_1 \rightarrow D_2$  be a function; unless otherwise stated  $f^* \in D_1^* \rightarrow D_2^*$  is defined as

$$f^*(\langle \rangle) = \langle \rangle, \quad f^*(d_1, \dots, d_n) = f(d_1), f(d_2), \dots, f(d_n)$$

### 3. Induction proofs

The basic proof principle employed for Scott's computable functions is that of fixed point induction. In this section we describe fixed point induction and its relation to other induction principles, structural and recursion induction.

#### 3.1. Fixed point induction

Fixed point induction [Pa69] appears in two disguises.<sup>5)</sup> First, suppose we are given an object which is defined as the least fixed point of a function,

<sup>5.)</sup> Actually, both instances of fixed point induction are identical since recursively defined domains are images under recursively defined retracts [Sc76b].

say  $f = \text{fix } h$ . If we want to prove that a property  $P$  holds for  $f$  we can do this by proving  $P(\perp)$  and  $\forall z.P(z) \Rightarrow P(h z)$ .

The second case obtains if we are given a recursively defined domain, say  $D = \Psi(D)$ , and wish to prove that property  $P$  holds for all elements in  $D$ . One can induct as follows:

- prove  $P(\perp)$
- assuming  $\forall z \in D.P(z)$  prove  $\forall z \in \Psi(D).P(z)$ .

In both cases it is necessary that property  $P$  is "admissible".

**[admissible]**  $P$  is admissible if and only if  $P(\bigcup f_i)$  whenever  $P(f_i)$  for all directed sets  $\{f_i\}$ .

To demonstrate the effect of admissibility consider  $g = \text{fix}(\lambda z.a z)$ , an infinite list of  $a$ 's. Consider the non-admissible proposition  $P \equiv \lambda z.z \neq a z$ . Clearly,  $P(\perp)$ . Assuming  $P(z)$ , i.e.  $z \neq a z$  it follows  $a z \neq a a z$  thus we have  $P(a z)$ . But observe that  $P$  is not true for the fixed point  $g$  since by definition of the fixed point we have  $g = a g$ .

### 3.2. Recursion induction, structural induction

We now show, that both these induction principles are special cases of fixed point induction.

Let  $f$  be a recursively defined function. Recursion induction proves that a proposition  $P$  holds for all results of  $f$ , given that  $f$  terminates. One proceeds by proving  $P$  for non recurring calls of  $f$  (base case) and then proving  $P$  for arbitrary calls to  $f$  assuming that inner calls satisfy  $P$ .

The basic difference to fixed point induction is, that recursion induction proves  $P$  only for terminating calls to  $f$ . Fixed point induction is applicable to nonterminating programs as well. In Scott's theory the result of an infinite computation is not necessarily undefined; rather is the limit point of an infinite sequence of partially computed values [Sc72b]. Since limit points in form of nonterminating computations are not considered in recursion induction, admissibility is not required.

Structural induction corresponds to fixed point induction on data structures. Again, the difference to fixed point induction is, that structural induction does not deal with infinite objects. It assumes that the domain in question is "well founded" [Ma74]. A domain is well founded if it has no infinite descending chains. For example, the domain of all finite lists is well founded, while the domain of all finite and infinite lists is not. In general,  $D = \Psi(D)$

is well founded if  $\Psi$  is strict. Since no infinite objects are present in a well founded set, admissibility is not required for structural induction.

Our definition of source and target languages contains both domains that admit structural induction and those that don't. The domains of the abstract syntax as well as modes allow structural induction, i.e. do not require admissibility. The domain of types<sup>6</sup> contains infinite chains in form of recursively defined types and structural induction is not applicable.

### 4. Verification techniques

We assume the reader to be familiar with the principles of Hoare's logic [Ho69] as well as the use of automated verification systems, such as the Stanford verifier [IL75, SV79].

This section deals with three problem areas which are not too well understood yet. We introduce new techniques to deal with pointers, quantification, and axiomatization of Scott's theory.

In general reasoning about pointers is extremely tedious. To isolate special simple cases we classify pointer operations as "reading", "updating", and "extension" operations; we prove a theorem that greatly simplifies reasoning about extension operations.

The second problem is that of expressing quantified formulas in a quantifier free assertion language. We present a deduction step in Hoare's logic called "weak  $\forall$ -introduction", comparable to  $\forall$ -introduction in first order logic.

Thirdly, we discuss the axiomatization of Scott's computable functions. This requires a satisfactory treatment of the undefined element ( $\perp$ ) as well as higher order functions.

#### 4.1. Pointers

##### 4.1.1. Pointers in the Stanford verifier

A satisfactory formal treatment of pointers has first been proposed by Luckham and Suzuki in [LS76]. The basic idea is to associate a "reference class" with each pointer type. The reference class is the collection of data objects pointed to by pointers of this type. The meaning of data structures involving pointers is explained relative to a reference class. The operation of dereferencing is described as "indexing" the reference class with a pointer. In 6.) abstract syntax, types, and modes are discussed in chapter III.

fact, the relation of pointers and reference classes is very much like that of indices and arrays.

Formally the expression  $p \uparrow$  means  $r \subset p \supset$  for the reference class  $r$  associated with the type of  $p$ . Assignments to data objects pointed to by pointers, e.g.  $p \uparrow \leftarrow e$ , are described as assignments to a whole reference class, e.g.  $r \leftarrow r, r \subset p \supset, e \supset$ . Similar to the array term  $\langle a, [i], e \rangle$  the term  $\langle r, r \subset p \supset, e \rangle$  denotes the reference class  $r$  after  $r \subset p \supset$  has been replaced by  $e$ . Assignment and selection obey the well known axioms

$$\begin{aligned} & \langle r, r \subset p \supset, e \rangle \supset C q \supset = \text{if } p = q \text{ then } e \text{ else } r \subset q \supset \\ & \quad \langle r, r \subset p \supset, r \subset p \supset \rangle = r \end{aligned}$$

Reference classes are not part of Pascal; the Stanford verifier introduces a reference class  $\#t$  automatically whenever the program contains a type declaration of the form  $pt = \uparrow t$ . In later language designs reference classes have been made explicit, e.g. Euclid [Lo78].

In addition to assignment and selection the "extension" operation  $r \cup p$  is defined for reference classes. This construct is used to describe the effect of generating a new data object; it extends an existing reference class  $r$  by the new cell  $p \uparrow$  pointed to by  $p$ .

#### 4.1.2. Reasoning about pointers

Suppose we are given a reference class  $r$  and a pointer  $p$  and we know that  $p$  points to the beginning of a linked list. If we now add a new element to this list for example, then this will change the reference class  $r$ . A central problem in reasoning about reference classes is to prove that this change does not destroy the already existing list structure. Proof techniques for this case have been outlined in [LS75]. However, these techniques are fairly tedious and the resulting complexity is unacceptable for our compiler proof.

We need more abstract operations on pointers that are easier to deal with. An obvious idea is to use the concept of data abstraction [Gu75, GH76] to encapsulate pointer operations and define modules which provide abstract list or tree operations. Still, the implementation of such a module has to proceed on a very low level and poses all of the above problems.

Another approach has been proposed recently [Su80]. Suzuki introduces the idea of pointer rotations. However, this method is not complete in the sense that all pointer operations could be expressed as rotation. Also, proof techniques for dealing with pointer rotations have yet to be developed.

Our approach to pointers is as follows. We distinguish three different operations on pointer structures, "reading", "updating", and "extending".

Special proof techniques for these three cases are described below.

"Reading" data stored in a pointer structure (e.g. linked list) causes no problem whatsoever, no change of reference classes occurs and thus all structures are preserved.

"Updating" on the other hand is extremely difficult to handle. An update operation is one which changes an existing reference class. The situation is familiar to the lisp programmer using language features like RPLACA and RPLACD. Update operations are only used at one place throughout the compiler proof. For this special case we prove that the particular update operation performed corresponds to the computation of a least fixed points. This situation is discussed in detail in section 6.

We talk about "extending" a pointer structure when all changes are of the form:

- create a new cell in the reference class.
- change this new cell.

Reasoning about extension operations is substantially simpler than that about update operations. This situation is analysed in the following subsection.

#### 4.1.3. Reasoning about extensions

Let  $pt = \uparrow t$  be a pointer type; as mentioned above,  $\#t$  is the associated reference class. Given a procedure  $extend$  which changes  $\#t$  by extension, then the verifier's semantic checking requires  $\#t$  to be declared *var global* to  $extend$  as in

```
procedure extend(pt:p);
global(var #t);
entry ...
```

If  $extend$  does not change any cells in the initial collection  $\#t_0$ , we would like any proposition  $P$  that was true for  $\#t_0$  to still be true for the extended reference class  $\#t_f$ . But of course, this is not correct. For instance, if  $P$  is a proposition about the cardinality of  $\#t$  or if it asserts that certain elements are not members of  $\#t$ .

We have to suitably limit the propositions  $P$  we allow to make the above statement true. Consider a pointer  $p$  which points to a linked list. Then we can define a function  $listrep(p, \#t)$  that maps  $p$  and  $\#t$  into an abstract list (the concept of representation functions is discussed in section 5 in greater detail).

If we have a proposition  $P$  not on  $\#t$  but on objects represented in terms of  $\#t$ , then we have the theorem that

$$P(\text{istrep}(p, \#t_0)) \Rightarrow P(\text{istrep}(p, \#t_f))$$

because  $\text{istrep}(p, \#t_0)$  must be independent of those cells which are not yet created and not part of  $\#t_0$ . This means that all objects defined in terms of the original reference class  $\#t_0$  are unchanged after *extend* has been executed. Consequently, if we disallow propositions on reference classes and only allow propositions on abstract objects defined in terms of reference classes the truth of any of these propositions is not effected by *extend*. This allows us to consider the reference class  $\#t$  a constant rather than a variable.

This situation can also be explained as follows. Given a program  $pgm$  that manipulates  $\#t$  by pure extensions. Suppose we had an oracle that at the beginning of the execution of  $pgm$  could guess the final reference class  $\#t_f$ . In this case we could initialize  $\#t$  to  $\#t_f$  and no extension operation would ever have to change  $\#t$ , i.e.  $\#t$  could be considered constant throughout  $pgm$ .

We can make use of the above observation in either of two ways. We can assume a meta theorem that allows to assume reference classes to be constant if they are only changes by extension operations. Alternatively, we can introduce suitable assertions that allow us to use the verifier to deduce relevant instances of this theorem.

We demonstrate the latter alternative in the proof of the parser in chapter IV. The above concepts are easily expressed in our assertion language using predicates *subclass* and *proper\_subclass*( $\#t_1, \#t_2$ ) is true, if  $\#t_2$  is an extension of  $\#t_1$ . *proper* is a predicate on abstract objects defined in terms of reference classes. It is true if such an object is well defined, i.e. if the representation function does not dereference an undefined pointer  $\#t \subset p \supset$  where the cell pointed to by  $p$  has not been added to  $\#t$ .

In the verification of the semantic analysis we do not repeat this reasoning, rather we resort to the meta theorem that allows us to consider a reference class unchanged by a procedure if the new reference class is a mere extension of the original reference class.

This reasoning simplifies proofs considerably and lets us concentrate on more relevant things.

#### 4.2. Quantification

In some cases it turns out that the quantifier free assertion language accepted by the Stanford verifier is a severe limitation. There are two principal methods for dealing with quantification.

A simple example illustrates this situation. Let us assume that we have a procedure  $p(y)$  for which we need the exit condition  $\forall z. R(z, y)$ . One solution is to introduce a new predicate symbol  $RQ(y)$  with the interpretation  $RQ(y) \equiv \forall z. R(z, y)$ . This method is general and allows elimination of all existential and universal quantifiers in arbitrary situations. The main disadvantage is, that all theorems about  $R$  that are only provable in a full first order theory have to be proven outside the verifier and supplied as lemmas.

In many cases a more elegant solution is possible which allows much of the proof to be handled by the verifier. Suppose that for a particular call to  $p$  we know that only one certain instance of the exit condition is used for the proof of this call, i.e.  $\forall z. R(z, y)$  is particularized to  $R(a, y)$  in the proof of this call. In this case we can write a procedure  $p(z, y)$  with exit condition  $R(z, y)$ . The additional virtual parameter  $z$  is used to provide the correct instance of the exit condition, e.g. if a particular call requires the instance  $R(a, y)$  of the exit condition for its proof, then we supply the virtual parameter  $a$  in this call. Note, that for this method no new predicates and lemmas are necessary.

This second method does not immediately apply for cases where we actually need the full quantification in the exit assertion. However, the following theorem allows us to verify a procedure in quantifier free logic using virtual parameters and then introduce quantifiers later on. This process is the analogue of the well known  $\vee$ -introduction in first order logic applied to Hoare's weak program logic.

Suppose we are given a procedure  $p(x_1, \dots, x_n, y_1, \dots, y_m)$  where all  $y_i$  are virtual parameters. Suppose further, that we have a proof for the body  $S$  of  $p$  with the specifications

$$\{Q(x_1, \dots, x_n, y_1, \dots, y_m)\}S(R(x_1, \dots, x_n, y_1, \dots, y_m))$$

Clearly, the final values of the var parameters among the  $x_i$  do not depend on the initial values of the virtual parameters  $y_i$ . Now we can construct a procedure  $p(x_1, \dots, x_n)$  with body  $S$  in which we omit all virtual code.

**Theorem** [weak  $\vee$ -introduction] Whenever

$$\{Q(x_1, \dots, x_n, y_1, \dots, y_m)\}S(R(x_1, \dots, x_n, y_1, \dots, y_m))$$

then it is true that

$$\{x_i = z_i\}S(\forall y_1, Q(z_1, \dots, z_n, y_1, \dots, y_m)) \Rightarrow R(z_1, \dots, z_n, y_1, \dots, y_m))$$

**Proof** Suppose the theorem were false, i.e. let  $z_i$  be such that  $S$  terminates

and assume that there is a tuple  $v_i$  with

$$Q(\bar{z}_1, \dots, \bar{z}_n, v_1, \dots, v_m) \wedge \neg R(z_1, \dots, z_n, v_1, \dots, v_m).$$

Clearly, the entry condition for the original procedure body  $S$  is satisfied and  $S$  will also terminate for  $\bar{z}$ , because termination must be independant of the values of virtual variables. But now by

$$Q(z_1, \dots, z_n, v_1, \dots, v_m) \{S\} R(z_1, \dots, z_n, v_1, \dots, v_m)$$

$R(z_1, \dots, z_n, v_1, \dots, v_m)$  must be true which contradicts our assumption. ■

Of course, we have to be able to express the new exit condition

$$\forall v_i. Q(\bar{z}_1, \dots, \bar{z}_n, v_1, \dots, v_m) \Rightarrow R(z_1, \dots, z_n, v_1, \dots, v_m)$$

in the assertion language. But this can be done by using the first method and introducing a new predicate symbol, say  $QR(\bar{z}_1, \dots, \bar{z}_n, z_1, \dots, z_n)$ . We still retain the advantage that for the proof of this exit condition we only need the definitions and properties of  $Q$  and  $R$  and that this proof is done in a quantifier free logic. To make any use of the new exit assertion  $QR$ , however, we need the theorem

$$\begin{aligned} QR(\bar{z}_1, \dots, \bar{z}_n, z_1, \dots, z_n) \wedge \\ Q(\bar{z}_1, \dots, \bar{z}_n, v_1, \dots, v_m) \end{aligned}$$

$$\Rightarrow R(z_1, \dots, z_n, v_1, \dots, v_m)$$

Application of this method is shown in chapter IV when we discuss code generation.

#### 4.3. Computable functions and first order logic

To specify the correctness of our compiler it is necessary to write assertions in first-order logic about functions over domains (spo's) of Scott's logic. The following section deals with the question how one can axiomatize such functions and prove useful properties about them.

##### 4.3.1. Standard interpretations

An obvious observation is that the usual axiomatizations of well known theories (e.g. Presburger arithmetic) do not allow corresponding functions from Scott's logic as models (e.g. addition over the flat domain of integers). For example, the sentence  $\forall z. z \neq z + 1$  is true for conventional addition but false for the flat integer CPO  $N_\perp$  and  $+ \in N_\perp \rightarrow N_\perp \rightarrow N_\perp$ .

This problem is relevant for our work since the verifier's prover has built-in decision procedures for certain theories, such as presburger arithmetic, lists, and data structures. Thus the system has + "reserved" as conventional addition and we have to introduce a new symbol to denote addition over the flat integer cpo.

In program proofs we systematically distinguish concrete types in the programming language from abstract domains in the compiler definition.<sup>7</sup> The relation between abstract and concrete is established through representation functions mapping objects of the program into abstract objects in the definitional language.

For instance, addition on  $N_\perp$  can be defined axiomatically as  $\text{plus}(x, y)$ . Using  $\text{plus}$  as addition on integers as it is known to the prover, a possible axiom set is:

$$\begin{aligned} \text{plus}(x, \perp) &= \perp \\ \text{plus}(\perp, y) &= \perp \\ \text{plus}(\text{intrep}(x), \text{intrep}(y)) &= \text{intrep}(x + y) \end{aligned}$$

where  $\text{intrep}$  is a representation function mapping integers into the integer cpo. With these axioms the equation

$$x = \text{plus}(x, \text{intrep}(1))$$

does not lead to inconsistencies. Rather it only allows to deduce

$$\neg \exists y. x = \text{intrep}(y).$$

#### 4.3.2. Higher order functions

In Scott's logic we have to deal with arbitrary high order function (in fact, elements of recursively defined domains can be considered "infinite" order functions [Sc72a]). The question is, how we can describe these objects in first order logic. The straightforward solution is to simply consider functions as values. All operations on functions, including function application have to be introduced as functions in the logic.

In addition we have to be able to talk about least fixed points of these functions. One solution is to axiomatize the ordering "less defined than" ( $\sqsubseteq$ ). Given this ordering it can be defined what it means to be the least fixed point etc. The main problem with this approach is, that things become excessively messy and incomprehensible and we should look for a simpler solution.

<sup>7</sup>) More on this in section 5

In [CM79] Cartwright and McCarthy propose a minimization schema to capture the idea of the least fixed point. But this approach is limited to first order functions and by introducing the predicate  $\text{id}_D$  (to test whether an element is proper) implicitly introduces Scott's ordering ( $\sqsubseteq$ ) for the special case of flat domains.

Here we go a different way. Given a denotational definition with a fixed set of domains. We say that a domain  $D$  is terminal if there is no other domain  $D'$  defined in terms of  $D$ . In particular, elements from a terminal domain can never appear as arguments to functions (otherwise the function would be an element of a domain defined in terms of this terminal domain).

For a given definition we let  $U$  be the sum of all non-terminal domains. All objects in  $U$  are axiomatized as values ( $U$  is the universe of the standard interpretation). Terminal functions from  $U \rightarrow U$  are first order functions. By the definition of terminal domains there are no higher order functions.

If functions are axiomatized as values we need an operation denoting function application. We use the symbol  $\otimes$  for this purpose; i.e. the standard interpretation for  $\otimes$  is  $T[[x \otimes y]\phi] = (T[[z]\phi)(T[[y]\phi])$ .

With these tools we can axiomatize non-recursive function definitions in a straightforward way. For example suppose we have a first order term  $\dot{T}(x)$  representing the term  $T(x)$  in Scottery. A definition  $g = \lambda z.T(z)$  is translated into the axiom  $g \otimes x = \dot{T}(x)$ .

#### 4.3.3. Least fixed points

For the above axiomatization of Scottery to be of any use we have to be able to axiomatize recursively defined functions and least fixed points. Suppose we are given  $f = \text{fix}(\lambda z.T(z))$  and  $\dot{T}$  as above. For  $f$  we can write the axiom  $\dot{T}(f) = f$  which is clearly satisfied by  $\text{fix}(\lambda z.T(z))$ . But the axiom merely asserts that  $f$  is a fixed point. It does not require that  $f$  be minimal; any fixed point satisfied the above axiom. We say that  $f$  is weakly axiomatized.

In many cases a weak axiomatization suffices to prove relevant properties of  $f$ . Suppose we can give a proof in first order logic of the form

$$f = g \otimes f \vdash P(f).$$

This means that  $P$  is true for all interpretations of  $f$  satisfying  $f = g \otimes f$ . In particular  $P(f \otimes g)$  is true. Although weak axiomatization is a useful device it is not sufficient in general and many properties  $P$  cannot be proven. Since  $\text{fix} \in (D \rightarrow D) \rightarrow D$  is a function we have the obvious property

$$f = g \Rightarrow \text{fix } f = \text{fix } g.$$

Again, a useful way of proving properties of fixed points. In section 6 we discuss yet another way to reason about fixed points. In certain cases it is possible to prove that a program computes a least fixed point without actually introducing the less defined relation  $\sqsubseteq$ .

#### 5. Representations

In this section we investigate how abstract objects can be represented by suitable data structures in a program. Clearly, in general the problem of data structure selection is very complex and requires human interaction. We are specifically interested in the representation of objects of Scott's computable functions. In this restricted domain several general rules are possible and will be discussed below.

A data structure is either a Pascal type or the product of types, written  $T_1 \times \dots \times T_n$ . Since such a tuple could always be defined as a record type we subsequently use the words type and data structure synonymously. We say that a domain  $D$  is represented by a type  $T$  by means of  $r$ , if  $r \in T \rightarrow D$  is a "representation" function mapping concrete objects of  $T$  into abstract objects in  $D$ . To emphasize that such a function has the special use as "representation function" we write  $r \in T \rightsquigarrow D$ . Semantically  $\rightsquigarrow$  is equivalent with  $\rightarrow$ ; syntactically  $\rightsquigarrow$  has lower precedence than  $\rightarrow$ . A representation function has to be continuous but can otherwise be arbitrary, we require neither surjectivity nor injectivity. Thus, for  $d \in D$  there may be none, one or several  $t \in T$  such that  $r t = d$ .

Unfortunately, our data types in Pascal are no domains; therefore, whenever we talk about functions defined on data types it is understood that the data type is a flat domain, i.e. has a bottom element added on. But we never use this bottom element to represent an element in a domain since this would not be an "effective" representation.

We now give several examples of types representing domains. More complicated representations will be introduced in the individual program proofs. The representations presented here are very natural ones and will be referred to later without being defined in each case.

#### 5.1. Primitive types

A scalar type  $T$  represents the flat domain  $T_\perp$  (up to isomorphism). The representation function  $\text{Rep} \in T \rightarrow T_\perp$  is a simple embedding.

For example, we have  $\text{truthrep} \in \text{boolean} \wedge^*(TT, FF) \perp$  with

$$\text{truthrep}(\text{true}) = TT, \quad \text{truthrep}(\text{false}) = FF.$$

A pointer type  $P = \uparrow T$  is the set of infinitely many pointers  $p$ , together with the element  $NIL$ . Let  $D_P = \{p_1, \dots, p_n, \dots\} \perp$ ; we have  $\text{rep}_P \in P \wedge^* D_P$  where  $\text{rep}_P = \lambda p. \text{If } p = NIL \text{ then } \perp \text{ else } p$ . In this case we have a representation of  $\perp$ .

### 5.2. Complex types

Let  $\text{rep}_1 \in T_1 \wedge^* D_1$  and  $T = \text{record } I_1:T_1; \dots; I_n:T_n \text{ end}$ . We define  $\text{rep} \in T \wedge^* D$  where  $D = D_1 \times \dots \times D_n$  as

$$\text{rep}(x) = (\text{rep}_1(x.I_1), \dots, \text{rep}_n(x.I_n)).$$

Had we union types in the verifier's language these would represent sums of domains in the very same manner. Since we do not have union types sum domains are represented as follows. Let  $\text{rep}_i \in T_i \wedge^* D_i$  and  $D = D_1 + \dots + D_n$ . With

$$T = \text{record tag:(1..n); } I_1.T_1; \dots; I_n.T_n \text{ end}$$

we define  $\text{rep} \in T \wedge^* D$  as  $\text{rep} = \lambda z. \text{let } i = z.\text{tag} \text{ in } \text{rep}_i(z.I_i).D$ .

Let  $T$  be a scalar type and  $\text{rep} \in T \wedge^* T \perp$  as above, then there is a function  $\text{rep}^{-1}$  such that  $\text{proper } \epsilon \equiv >\text{rep}(\text{rep}^{-1}\epsilon) \equiv \epsilon$ .

Let  $\text{rep}_1 \in T_1 \wedge^* D_1$  and  $\text{rep}_2 \in T_2 \wedge^* D_2$ . With  $T = \text{array}[T_1] \text{ of } T_2$  and  $D = D_1 \rightarrow D_2$  we can define a representation  $\text{rep} \in T \wedge^* D$ . By Pascal rules  $T_1$  must be a scalar type, thus the following is well defined:

$$\text{rep}(a) = \lambda \epsilon \in D_1. \text{ If } \text{proper } \epsilon \text{ then } \text{rep}_2(a[\text{rep}_1^{-1}(\epsilon)]) \text{ else } \perp.$$

To make any sense, of course,  $\text{rep}(a)$  has to be continuous for arbitrary arrays  $a$ . This is in fact guaranteed since  $\text{rep}(a)$  is strict and defined on a flat domain.

Declaring  $\uparrow T$ , the verifier automatically introduces a type  $\#T$ , the "reference class" for  $T$ . As mentioned earlier a reference class can be thought of as an array indexed by pointers in  $\uparrow T$  with  $T$  as its component type. Given  $T$ ,  $\uparrow T$ ,  $\#T$  then by the above remarks we have

- $\text{rep} \in T \wedge^* D$ ,
- $\text{rep}_P \in \uparrow T \wedge^* D_P$ ,
- $\# \text{rep} \in \#T \wedge^* D_P \rightarrow D$ , and

In general, a domain need not be represented by a concrete data structure immediately. Rather, it suffices to have a representation for a domain in terms of other domains for which there are concrete representations.

### 5.3. Recursive domains

Recursively defined domains can also be represented in a straightforward way. Consider the simple example  $D = A + (B \times D)$ . By the above principle we can construct a representation  $T$  with  $\text{rep} \in T \wedge^* A + (B \times D_F)$  with  $\# \text{rep} \in \#T \wedge^* D_F \rightarrow D$ , and  $\text{rep}_P \in \#T \wedge^* A + (B \times D_F)$  with  $\# \text{rep} \in \#T \wedge^* D_F$ .

Now consider the data structure  $(\uparrow T, \#T)$  consisting of pairs of a pointer  $P \in \uparrow T$  and a reference class  $r \in \#T$ ;  $(\uparrow T, \#T)$  represents  $D$  by means of  $\text{rep} \in (\uparrow T \times \#T) \rightarrow D$  as follows.

$$\begin{aligned} \text{recrep}(p, r) = & \text{ let } z = (\# \text{rep } r)(\text{rep } p) \text{ in} \\ & \text{If } x \in A \text{ then } z:A \text{ else} \\ & \langle x^{\#1}, \text{recrep}(\text{rep}_P^{-1}z^{\#2}, r) \rangle \end{aligned}$$

### 6. Fixed points

#### 6.1. Reasoning about fixed points

At several places we give manual proofs about properties of fixed points. Sometimes it is more natural to reason about recursion equations than to argue about the fixed point operator ( $\text{fix}$ ). This is particularly true if "simultaneous" recursion is required. For this reason we introduce the following notation as an alternate for the fixed point operator.

A term of the form

$$\left( \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right) \leftarrow \left( \begin{array}{c} T_1(x_1, \dots, x_n) \\ T_2(x_1, \dots, x_n) \\ \vdots \\ T_n(x_1, \dots, x_n) \end{array} \right)$$

where  $T_i$  are arbitrary terms that may involve  $x_j$ 's is equivalent to

$$(\text{fix } \lambda x. (T_1(x^{\#1}, \dots, x^{\#n}), \dots, T_n(x^{\#1}, \dots, x^{\#n})))^{\#1}$$

The following syntactic operations can easily be shown to be sound.

- Row introduction -- If terms  $T_i$  contain a term  $T$  then several occurrences of  $T$  can be replaced by a new variable  $x_{n+1}$  if the new row  $x_{n+1} \leftarrow T$  is added to the system.<sup>8</sup>
- observing some obvious restrictions concerning free and bound variables

- Elimination of redundant rows — If  $x_i, i \neq 1$  does not occur in any term except possibly in  $T$ , then row  $i$  can be deleted from the system.
- Elimination of duplicate rows — If two rows  $i$  and  $j$ ,  $i, j \neq 1$  are identical up to a renaming of the variables  $x_i$  and  $x_j$ , then one of these rows, say  $i$ , can be deleted if all occurrences of  $x_i$  are replaced by  $x_j$  in the remaining system.

• Substitution — Given a row  $z \leftarrow T$ , then any number of occurrences of  $z$  in any  $T$ , may be replaced by  $T$ .

We use these rules to prove equalities of the form  $\text{fix } f = \text{fix } g$  by rewriting both fixed points as recursion equations and transforming them into isomorphic equations. Although this works in many cases it is by no means complete; in general fixed point induction is required.

### 6.2. Operationalization of fixed points

In the compiler to be written it is necessary to compute least fixed points. For example, in LS one can define a type which is a pointer to a type which is a pointer to a type which is a pointer ... etc. The semantics of such a recursive type is given as a least fixed point. In order to do typechecking in the compiler we have to compute a representation of this least fixed point. The weak axiomatization introduced in 4.3 is insufficient to prove that a particular computation computes the "least" fixed point. The representation theory presented above allows a way out. Instead of axiomatizing a least fixed point we axiomatize a representation of this fixed point. It turns out that in many cases this representation can be described without the use of  $\text{fix}$ .

As a simple example consider the domain  $D = A + (B \times D)$  for which we defined a representation  $\text{recrep}$  in the previous section. Let  $z \in T$  be such that  $\text{rep}(z) = (a, y)$  and  $\text{rep}^{-1}(y) = p$ . With the assignment  $p \mapsto z$  we can compute a reference class  $r = < r, \subset p \supset, z >$ . Applying  $\text{recrep}$  we get

$$\text{recrep}(p, r) = \text{fix}(\lambda y. (a, y)).$$

Thus the above theory gives us the possibility to compute the representation of least fixed points without using the least fixed point operator  $\text{fix}$  directly.

We call this process "operationalization" of fixed points. A similar idea more related to reasoning about programs has been proposed by M. Gordon in [Go'75].

The above example may appear to be a special situation; but this is not so, we now prove a more general theorem showing that operationalization can be employed systematically in a variety of cases. Further research has to

determine what the most general situation is that allows for operationalization but the theorem presented below suffices for all cases arising in our work.

In general we proceed as follows: for a domain  $D$  to define a representation  $r \in T \wedge^* D$  and a representation  $f \in \hat{T} \wedge^* D \rightarrow D$ . We then prove that there is an operation  $\Phi$  on  $\hat{T}$ , producing an object in  $T$  such that  $r(\Phi f) = \text{fix}(f)$ ; i.e. the following diagram commutes:

$$\begin{array}{ccc} & \Phi & \\ \text{---} & \nearrow & \searrow \\ r & & \text{fix} \\ \text{---} & \searrow & \nearrow \\ & f & \end{array}$$

1. let  $D$  be a recursive domain. We say that a family of function  $g_i \in D^n \rightarrow D$  generates  $D$ , if every element of  $D$  is given by some composition of  $g_i$  and  $\perp$  (some of the  $g_i$  may be constant). This situation obtains in our compiler for many domains, e.g. abstract syntax, modes, and types. Of course,  $g_i$  need not all be of the same arity, also, they may take additional parameters from domains not depending on  $D$ . What follows adapts easily to either situation.

$D$  can be represented as outlined in the previous section. Let  $T = N \times P \times \dots \times P$  where  $P = \# T$ ; let  $\# T$  be the reference class for  $P$ . Now  $D$  is represented by pairs in  $\# T \times P$  by means of  $r \in \# T \rightarrow P \wedge^* D$  as

$$\begin{aligned} r \circ p = & (\text{fix } \lambda h. \lambda x. \text{let } (i, p_1, \dots, p_n) = c \circ x \text{ in} \\ & \quad \text{let } y_j = h(p_j \text{ in } g(y_1, \dots, y_n))p \\ & \quad \text{in } g(y_1, \dots, y_n))p \end{aligned}$$

Furthermore we can define a representation for the domain  $D \rightarrow D$  in terms of triples  $\# T \times P \times P$  with  $f \in \# T \rightarrow P \wedge^* D \rightarrow D$  as

$$\begin{aligned} f \circ p q = & (\text{fix } \lambda f. \lambda x. \lambda d. \text{let } (i, p_1, \dots, p_n) = c \circ x \text{ in} \\ & \quad \text{let } y_j = \text{if } z = q \text{ then } c \circ p \text{ else } f(p_j, d) \\ & \quad \text{in } g(y_1, \dots, y_n))p \end{aligned}$$

Clearly, not all functions in  $D \rightarrow D$  can be represented in this way but many useful ones can (recall that we do not require representations to be surjective). We now have the important

**Theorem** For  $r$  and  $f$  as defined above:

$$r(c[c/p/q])p = \text{fix}(f \circ p q)$$

**Proof** We express the least fixed points as recursion equations as introduced in section 6.1. Observing

$$\begin{aligned} r(c[p/q]) p &= (\text{fix } \lambda h. \lambda z. \\ &\quad \text{let } (i, p_1, \dots, p_n) = \text{if } z = q \text{ then } c p \text{ else } c x \text{ in} \\ &\quad \text{let } y_j = h p_j \text{ in } g(y_1, \dots, y_n)) p \end{aligned}$$

the left hand side of the theorem becomes

$$\binom{z_1}{h} \leftarrow \binom{h \ p}{\lambda z. \text{let } (i, p_1, \dots, p_n) = \text{if } z = q \text{ then } c p \text{ else } c x \text{ in} \\ \text{let } y_j = h p_j \text{ in } g(y_1, \dots, y_n)}$$

The right hand side may be rewritten as

$$\binom{z_2}{f} \leftarrow \binom{k \ p}{\lambda z. \text{let } (i, p_1, \dots, p_n) = c z \text{ in} \\ \text{let } y_j = \text{if } p_j = q \text{ then } z_2 \text{ else } f \ p_j \ z_2 \text{ in } g(y_1, \dots, y_n)}$$

Let us introduce  $k = \lambda p. f \ p \ z_1$  in this system yielding

$$\binom{z_2}{k} \leftarrow \binom{k \ p}{\lambda z. \text{let } (i, p_1, \dots, p_n) = c z \text{ in} \\ \text{let } y_j = \text{if } p_j = q \text{ then } k \ p_j \text{ in } g(y_1, \dots, y_n)}$$

Now let  $t \ p_j = \text{if } p_j = q \text{ then } k \ p_j$  in this system, thus

$$\binom{z_2}{t} \leftarrow \binom{k \ p}{\lambda z. \text{let } (i, p_1, \dots, p_n) = c p \text{ in} \\ \text{let } y_j = t \ p_j \text{ in } g(y_1, \dots, y_n) \\ \text{let } y_j = t \ p_j \text{ in } g(y_1, \dots, y_n)}$$

which proves the equality.

## 7. Intermediate representation of programs

In our compiler the source program is transformed into various intermediate representations. These are token sequences, syntax trees, and abstract syntax. The concepts of sequences is well understood; this section deals with trees and abstract syntax.

## 7.2. Abstract syntax

The idea of an abstract syntax has first been introduced by McCarthy in [Mc66]. It allows to reason about programs regardless of their external representation (on paper or in a machine). Formally, abstract syntax can be defined as a free algebra, for our purposes abstract syntax is a domain  $\sigma$  of programs.

9.) This pun is intentional: the set of  $\Sigma$ -trees on  $X$  is in fact the  $\Sigma$ -word algebra on  $X$ ; see [Cof65].

## 7.1. Trees

### 7.1.1. $\Sigma$ -trees on $X$

[Trees] Let  $\Sigma$  be a finite alphabet and  $\sigma \in \Sigma \mapsto N$  a function assigning an “arity” to each  $\sigma \in \Sigma$ . Let  $X$  be an arbitrary set, then  $W_\Sigma(X)$ , the set of  $\Sigma$ -trees on  $X^0$  is defined as follows:

- $X \subseteq W_\Sigma(X)$
  - if  $r_1, \dots, r_n \in W_\Sigma(X)$  and  $\sigma(\sigma) = n$  then  $(\sigma r_1 \dots r_n)$  is in  $W_\Sigma(X)$
  - These are all elements in  $W_\Sigma(X)$
- If  $X$  is a domain, then an analogous construction can be devised such that  $W_\Sigma(X)$  is a domain too.
- 7.1.2. Operations on trees**
- leaves**  $\in W_\Sigma(X) \mapsto X^*$
- leaves  $x = \{\sigma\}$**
- leaves  $(\sigma r_1 \dots r_n) = \{\sigma\} \cup \text{leaves } r_1 \dots \cup \text{leaves } r_n$**
- root**  $\in W_\Sigma(X) \mapsto \Sigma \cup X$
- root  $x = x$**
- root  $(\sigma r_1 \dots r_n) = \sigma$**

The domain **asyn** is defined recursively using strict products. Therefore, it is a flat domain of finite trees. **asyn** is a subset of  $W_{\Sigma}(X)$  where  $X$  is a set of primitive concepts (like identifiers, numerals etc.) and  $\Sigma$  is a set of constructors to build more complex programs. (To characterize **asyn** more precisely we had to introduce many sorted algebras).

For example, if  $E$  and  $\Gamma$  are trees then  $(\text{while } E \ \Gamma)$  is an element (tree) in the domain of programs. To increase readability we use an informal notation resembling the external representation of trees. For example, we write

$\text{while } E \text{ do } \Gamma$

instead of  $(\text{while } E \ \Gamma)$ .

### 7.2.1. Constructor and selector functions

Let **zzz** be a name for an element in  $\Sigma$  with  $a(\text{zzz}) = n$ ; we agree to use the symbol **mkezzz** as a  $n$ -ary function from  $(W_{\Sigma}(X))^n$  to  $W_{\Sigma}(X)$  "constructing" a new tree according to

$$\text{mkezzz}(r_1, \dots, r_n) = (\text{zzz } r_1 \dots r_n).$$

Furthermore we agree that **deczzz** is the inverse operation to **mkezzz** which returns the list of subtrees of a given tree:

$$\text{deczzz}(\text{mkezzz}(r_1, \dots, r_n)) = r_1 \dots r_n.$$

Similarly we use the convention that **iszzz** is a predicate which is true for all trees with root **zzz** and false otherwise.

In later sections we intermix the constructor notation with the above informal one. For example, we use "while  $E$  do  $\Gamma$  od" in the formal language definition while we use " $\text{mkezzz}(E, \Gamma)$ " in the machine readable lemmas input to the verifier.

### 7.2.2. Conformity relations

The **conformity relation** :: is used to conveniently access subtrees of a tree. Given  $W_{\Sigma}(X)$  with **zzz** being a name for an element in  $\Sigma$  then

$$E :: \text{mkezzz}(x_1, \dots, x_n, y_1, \dots, y_m)$$

is true in a context where all  $y_i$  are bound and  $x_i$  are free if and only if  $\exists x_1 \dots x_n. E = \text{mkezzz}(x_1, \dots, x_n, y_1, \dots, y_m)$ . In addition, conformity relations cause free variables to be bound to corresponding components. For example,

$$\text{if } E_0 :: \text{mkezzz}(x_1, \dots, x_n, y_1, \dots, y_m) \text{ then } E_1 \text{ else } E_2$$

means

$\text{if iszzz}(E_0) \text{ then}$ $\quad \text{let } x_i = \text{deczzz}(E)^{\#i} \text{ in } E_1$ $\text{else } E_2$	$\text{where again } x_i \text{ are free variables.}$
---	---

### Chapter III. Source and target languages

In this chapter we introduce source and target languages *LS* and *LT*. We first give an informal overview over the source language followed by its formal definition. The target language is described in a similar fashion.

#### 1. The source language

One objective of this work is to write a compiler for a realistic and useful language. Since Pascal is a fairly simple but still useful language our source language *LS* is based on Pascal.

A second criterium is that we need a formal definition of *LS*. There exist formal definitions of Pascal of various flavours [HW73, AA77, Te77a]. We build on these, notably Tennent's denotational semantics [Te77a].

However, we made some changes to Pascal. Several language features of Pascal have been omitted. This has quantitative rather than qualitative reasons and was done to keep the task of verifying the compiler manageable. The theory and proof techniques developed in this thesis is general enough to handle all omitted cases.

Some parts of Pascal are poorly defined (see [Ha73, WS77]) and have been changed or extended. Some of the extensions follow naturally from an attempt to formally define the language as has been pointed out in [Te77b]. This section gives an informal overview of *LS* and points out differences to Pascal.

##### 1.1. Data structures

As Habermann [Ha73] points out, the notion of type is not very well defined in Pascal. Aside from his largely informal arguments this also becomes apparent if one tries to give a formal definition of the type concept of Pascal. In fact, in Tennent's definition the Pascal concept is completely abandoned and a structural type equality in the sense of Algol 68 is substituted.

Types in Algol 68 are equal if they are of the same structure. That is two records are equal if they have the same field selectors and equal selectors refer to components of equal type. Two pointer types are equal if they point to equal types.

This concept seems appealing, however, it poses some problems. It is non-trivial to determine the equality of two types. Engelfield [En72] has shown that this is equivalent to determine equivalent states of a finite automaton.<sup>1</sup>

Defining "structural" equality formally is not at all straightforward. In the original report [Wi69] type equality was "defined" in English while the revised report [Wi76] uses a non-intuitive definition in terms of van Wijngaarden predicates. In Tennent's denotational definition of Pascal [Te77a] the Algol 68 style equality is used. One problem in his semantics is that the definition of a monotonic equality on types requires to count the number of types defined in a program which seems rather unnatural. Consider the two types type  $P = \uparrow q; q = \uparrow p$ . Two types are equal in the Algol 68 sense if they are both pointers to equal types, thus we have  $\text{equal}(p, q) \equiv \text{equal}(q, p)$  which is not well defined. Tennent uses the number  $n$  of types defined in an environment and defines  $\text{equal}(p, q, n) = \text{if } n = 0 \text{ then true else } \text{equal}(q, p, n - 1)$ .

Pascal's types are very simple to compare; each type declaration defines a new unique type, not equal to any other type. Besides being easy to implement it allows the programmer to take advantage of additional type checking possibilities: one may have two types array [1..5] of integer which are different. The Pascal view of types causes several problems if we consider subranges. A constant of an enumerated type may also be of a subrange type. Operations on enumerated types must be viewed as generic (for all subranges) or elements of a subrange have to be coerced to the supertype before such operations can be applied. Habermann [Ha73] presents various examples involving these problems.

Our solution to types is as follows. Two type declarations define different types as in Pascal. However, subranges do not constitute new types; rather, a subrange type is equal to the supertype of which it is a subrange. Each scalar type in *LS* has a "subrange attribute" which specifies the allowable range of values of this type. Variables of type subrange are considered variables of the supertype except that assignments to them causes a runtime check for the bounds.

*LS* has the predefined types Integer and Boolean. It allows the definition of enumerated types and subranges [JW76]. Furthermore the user can define arrays.

1.) Actually the problem is further complicated by union types which are left out of Engelfield's considerations.

rays with constant bounds, records without variants, and pointer types. Even though variable bounds for arrays would be a reasonable extension and make the implementation slightly more interesting they have been omitted. Adding variable bounds means that subranges can no longer be used to define the index range and a new semantic concept has to be introduced in the language. Notable omissions from Pascal are characters, reals and file types. Also, we do not consider records with variants as their definition is Pascal is unsafe (i.e. violates strong typing).

*L* allows the definition of constants of type boolean and integer. This concept is a slight generalization of Pascal in that an arbitrary expression is allowed to define a constant provided this expression can be evaluated at compile time.

A major deficiency of Pascal is that there are no constant denotations for arrays and records. In fact, the type concept of Pascal makes it hard to incorporate constant arrays. For example what should the type be of say [1, 2, 3, 4, 5]. Certainly it is not array [1..5] of integer; one possible solution would require constructors for constants explicitly naming the type. Constant denotations for complex types are omitted in *L*.

### 1.2. Program structures

Statements of *L* are those of Pascal with the exception of for-loops, case and with-statements. These are omitted since they are not strictly necessary in the language and their implementation poses no particular problem (possible extensions of our language to include the above features are discussed in chapter V). The while and if statement are changed to contain a syntactical end symbol which allows while bodies and if branches to be statement lists without having to write begin and end.

Pascal's block structure has been extended to allow declarations in inner blocks. *L* allows jumps as in Pascal; jumps into blocks or statements are illegal. Labels do not have to be declared and may be reused in inner scopes.

Procedures and functions are recursive but not mutually recursive. The reason is that we omit forward declarations of procedures. Alternatively, we could allow procedure and function identifiers to be used before their declaration. Either extension would merely complicate the static semantics; our code generation can handle mutual calls.

Procedures and functions may not be passed as parameters. One problem with procedure parameters is that this mechanism in its full generality requires type checking at runtime. Alternatively one could require the specification of

the parameter types of the formal procedure parameters as in Algol68. For example

```
procedure p(var z:integer; f:procedure(integer,boolean);integer);
```

Each *L* program has one anonymous input and one anonymous output file both of type integer. *L* provides the operations read, eof, and write operating on these files. Opening and closing of files is done automatically at the beginning and end of program execution. Definition, implementation and proof can be extended easily for other predefined functions and procedures; more I/O operations can be added. An extension of the input output facilities to those of Pascal requires the introduction of file types as well as the formal definition of certain components of the operating system environment. For example the semantics of opening a nonexistent file has to be defined. Expressions in *L* are those of Pascal with exception of the precedence of operators which is more natural in *L*.

### 2. Formal definition of *L*

#### 2.1. Structure of the formal definition

As outlined earlier a language is defined in several stages: micro syntax, syntax, tree transformation, and semantics. Semantics in turn will be subdivided into static and dynamic semantics. Each of these components defines a mapping; the composition of these mappings defines the meaning of a program. We now define the individual mappings and their composition.

The following domains are common to all components of the compiler. Internally, the definitions of individual components will utilize additional domains.

$c \in C^h$	Characters
$s \in S^h = C^h$	Strings
$T^h$	Terminal symbols
$V^h$	Token values
$T^k = T^m \times V^l$	Tokens
$S^t$	Syntax trees
$P^m$	Abstract programs
$E_T$	Token errors
$E_S$	Syntax errors

$E_C$   
 $E_R$   
 $E = E_T + E_S + E_C + E_R$

Compile time errors  
 Runtime errors

$M = N^* \rightarrow (N^* + E_R)$   
 $M = N^*$

Meaning of programs  
 Most of these domains are explained in more detail in the corresponding section. For the present discussion their internal structure is irrelevant. It was necessary on this level to define  $M$ , the meaning of programs, since this shows that errors can not only occur during compilation but also during execution. The meaning of a program then is a mapping from the input file,  $N^*$ , into the output file  $N^*$  or a possible error.

The language is completely specified by defining the above domains and the following functions:

$scan \in Str \rightarrow (Tk^* + Err)$  Scanner  
 $parse \in Tk^* \rightarrow (St + E_S)$  Parser  
 $treeir \in St \rightarrow Pgm$  Tree transformations  
 $ssem \in Pgm \rightarrow PgM + E_C$  Static semantics  
 $dsem \in PgM \rightarrow M + E_R$  Dynamic semantics

[ $E$ -strict function composition] Let  $f \in D_2 \rightarrow D_3$  and  $g \in D_1 \rightarrow D_2 + E$  then  $E$ -strict function composition,  $f \odot g \in D_1 \rightarrow D_3 + E$  is defined as

$$f \odot g = \lambda x. \text{let } t = g x \text{ in if } t \in E \text{ then } t \text{ else } f t$$

The meaning of a program is a mapping from character strings into  $M$  given as:

$$smean \in Str \rightarrow M + E$$

$$smean = dsem \odot ssem \odot treeir \odot parse \odot scan$$

## 2.2. Denotational semantics

The definitional formalisms used to define micro syntax, syntax, and tree transformations are fairly well known. Denotational semantics, however, is not yet common knowledge. Therefore we review the basic principles at this point. Both static and dynamic semantics are defined using denotational semantics.

### 2.2.1. General concepts

A denotational language definition is a mapping which assigns a "meaning" to each program.<sup>2</sup> Here program means an abstract program in some abstract syntax as defined above.

The next question of course is, what is the meaning of a program. The method leaves ample room for different meanings. In general, the meaning of a program is a mapping from input data to a set of answers. Answers in turn can be of different flavour. For example, an answer could be an output file or a runtime error. Alternatively, answers could be *true* and *false* merely indicating termination or nondetermination.

Note, that a denotational definition is not an interpreter. The meaning of a program is a (function) object which can be reasoned about. The meaning of a program is *purely extensional*. This means that two programs are considered equal whenever they denote the same function. For example, questions as to the complexity of a program or order of evaluation cannot be answered on the basis of a denotational definition. This in turn has important advantages. In particular in the case of a compiler it leaves ample room for optimizations. As indicated in the introduction, we will prove that source and target program have the same meaning. Thus, sophisticated code optimization could be incorporated in our compiler without having to change the formal language definitions.

One part of a denotational definition is the definition of the underlying domains. By specifying the domains we define the relevant semantic concepts and their relationship. For example, such concepts as types, storable values, locations, answers etc. are specified.

A set of "valuation functions" assigns a meaning to language constructs in terms of the semantic concepts.

The comparison to first order logic may help clarify the situation (see chapter II, section 1.2). Formulas are syntactic objects. An interpretation assigns a meaning (a truthvalue in this case) to each formula. Semantic domains are fairly simple, we only have one universe. But this need not be so, just consider many sorted logical systems.

### 2.2.2. Semantic concepts of Algol-like languages

Imperative languages such as Pascal and Algol are an outgrowth of the vonNeumann architecture. Therefore their semantics is most naturally described in terms of *locations*, *values* stored in locations and a *memory*, here called

2.) A program denotes its meaning, hence the name denotational semantics.

**store.** Let  $L$  and  $V$  be the domains of locations and values respectively.<sup>3</sup> A store can then be defined as a mapping from  $L$  to  $V$ , i.e. we write  $S = L \rightarrow V$ ; it assigns a value to each location.

Another concept is that of identifiers used in a program. Identifiers name variables, procedures, functions and the like. An environment is used to keep track of all the current definitions of identifiers. Considering only variables for the time being, an environment is a mapping which assigns a location to each variable identifiers; we write  $U = Id \rightarrow L$ . In practice environments are more complicated.

What is the meaning (value) of a statement? A naive answer is, that a statement is a mapping from stores to stores; a statement maps the store before its execution into the store after its execution. However, this approach is inadequate to define the meaning of statements that do not normally terminate or alter the "normal" control flow with a jump. Consider the two statements  $\Gamma_1, \Gamma_2$ . We would like to describe the meaning of this sequence of two statements as the composition of their individual meaning. Assume that  $\Gamma_1$  and  $\Gamma_2$  are defined by a mappings  $f_1, f_2 \in S \rightarrow S$  respectively. The meaning of  $\Gamma_1, \Gamma_2$  would be the mapping  $f_2 \circ f_1$ . Suppose  $\Gamma_1$  does not terminate normally but instead produces a runtime error or a jump to another part of the program; our simple model clearly fails:  $f_2$  will not be applied to the result of  $f_1$  since control never reaches  $\Gamma_2$  in this case.

To remedy this situation continuations have been introduced [SW74]; they are based on an idea derived from "tailrecursion" [Ma71]. A continuation is a function from stores into answers, i.e.  $C = S \rightarrow A$ . At a given point during program execution a continuation specifies the behaviour of the remainder of the program. This behaviour is a mapping of the current store to the final result (answer) of the program. The way to describe the meaning of statements now is as follows. Suppose we know the continuation that obtains after a statement  $\Gamma$ , that is we know what the answer of the program will be once  $\Gamma$  terminated with a given store. Now we can determine the continuation that obtains before the execution of  $\Gamma$ . Thus the meaning of statements will be a mapping from continuations to continuations ( $C \rightarrow C$ ). But note, that the reasoning proceeds backwards, i.e. the meaning of a statement  $\Gamma$  maps the continuation "after"  $\Gamma$  into that "before"  $\Gamma$ .

Using continuations the description of runtime errors and jumps poses no problem at all. Consider the sequence  $x \leftarrow a / 0 ; \Gamma$ . The meaning of  $x \leftarrow a / 0$  will be  $\lambda \theta . \text{divide\_error}$ . That is, whatever continuation  $\theta$  holds before  $\Gamma$  and

3.) these are further specified in section 6.

after  $x \leftarrow a / 0$ , the continuation before  $x \leftarrow a / 0$  will be  $\lambda \sigma . \text{divide\_error}$ . I.e. no matter what state we are in when we execute  $x \leftarrow a / 0$ , the answer of the program is  $\text{divide\_error}$ .

A similar situation exists for expressions. A naive view is to consider an expression as a mapping from stores to values. Again, this would not capture side effects of expressions, runtime errors, and jumps out of expressions. The solution is the introduction of expression continuations  $K = V \rightarrow C$ . Intuitively, an expression continuation tells us what the final answer of the program will be given a value (the result of an expression) and a store; it tells us what "to do" with the expression value if we succeed in computing it. The meaning of an expression then is a mapping from expression continuations to continuations. Given the expression  $E$  and the expression continuation  $\kappa$  and suppose that  $E$  evaluates to  $\epsilon$ , then the resulting continuation will be  $\kappa\epsilon$ . On the other hand, if  $E$  causes a runtime error, the continuation will be  $\lambda \sigma . \text{error}$ . In a similar fashion more complex domains are constructed; examples are function and procedure values. For example, a procedure is a mapping in  $V^* \rightarrow C \rightarrow C$ , i.e. it takes a list of parameter values and a continuation and produces a new continuation.

Continuations, expression continuations, procedure, and function values all are domains with a similar structure. In the definition of LS we introduce the novel concept of generalized continuations which combines all of the above domains into one and leads to some simplification of the language definition (see section 7).

### 2.2.3. Denotational definition of a machine language

Concepts used to define high level languages can also be used to define machine languages. Since we have no variable names the environment is significantly simpler. In section 8 we describe a stack machine and define appropriate domains. There will be stacks  $S$ , memories  $M$ , and displays  $D$ . Continuations in the target language are mappings of the form  $D \rightarrow S \rightarrow M \rightarrow A$ . The intuitive meaning of such a continuation is similar to that in the source language case: given a display, a stack, and a memory, the continuation yields the final answer of the program.

### 2.2.4. Notational issues

A main criticism of denotational semantics is that the notation is too complicated and confusing. We have adopted the notation introduced initially by Scott and Strachey [SS71] and feel that it is very appropriate.

Since the formulas we are dealing with are extremely large, a short notation is imperative for comprehension and overview. Definitions and proofs are more transparent and smaller. As a trade off we have to learn and get familiar with the new notation.

The main points of our notation can be summarized as follows. Function application is written as juxtaposition as in  $f \cdot y$ . The semicolon is used to break the normal precedence of function application as in  $f \cdot x; g \cdot z$ . Special parentheses  $[ \dots ]$  are introduced to distinguish syntactic from semantics objects. Attempts are made to have variables named by single letters; we use of roman, greek, and script fonts in upper and lower case.

Mathematics has produced plenty of examples where new short notations have been fruitful. Calculus without  $\int$ ,  $\frac{df}{dx}$  is unthinkable.

Examples of denotational definitions written in different notations are not very readable, more error prone, and of less help in reasoning about the language [DK79, Pa79, BJ78].

### 3. Micro syntax

The micro syntax defines a mapping from a sequence of characters into a sequence of tokens. A token is a pair containing syntactic and semantic information.

#### 3.1. Definitional formalism

Micro syntax is defined in two steps. First we define which substrings of the input constitute single tokens. Each of these strings is then mapped into the corresponding syntax semantic pair.

We consider a finite number of classes of tokens; the external representation of elements of each class is given as a regular language  $L_i$ . For example an identifier is a string of digits and letters, starting with a letter. Similarly, numbers and single and multi-character delimiters are defined.

Let  $\{L_i\}$  be a set of regular languages each defining a class of tokens. We define rules according to which the input string is to be broken into substrings each of which is in one language  $L_i$ . The rule is to find the longest leftmost substring  $s$  of the input such that  $s \in L_i$ , for some  $i$ . Given  $s \in L_i$ , a semantic function  $S_i$  maps the string  $s$  into the token represented by  $s$ .

For example,  $L_{id}$  is the language of identifiers.  $S_{id}$  checks if an identifier is a reserved word and treats it differently from ordinary identifiers. Further, an

identifier which is not reserved is encoded as a unique integer by  $S_{id}$ . Clearly, to do this the functionality  $L_{id} \rightarrow Tk$  is too simple for  $S_{id}$ . Rather  $S_{id}$  needs additional information about identifiers that have been scanned previously as well as which codes have already been used to encode identifiers.

We introduce identifier tables

$$h \in H = (L_{id} \rightarrow N) \times (N \rightarrow \{\text{used}, \text{unused}\})$$

mapping identifiers into  $N$  and specify whether a number in  $N$  is used or unused. Now, the functionality of  $S_i$  becomes

$$S_i \in H \rightarrow L_i \rightarrow (Tk \times H).$$

Let us now describe how a given set of languages  $L_i$ , and semantics functions  $S_i$  are combined to define the scanning function  $scan \in Str \rightarrow (Tk^* \rightarrow Err)$ . We define

$$scan = \Psi h_0$$

where  $h_0$  is the initial identifier table which contains all predefined identifiers of the language.

$\Psi \in H \rightarrow Str \rightarrow Tk^*$  deletes initial substrings from the input string until the first character  $c$  is the beginning of a token.  $c$  is the beginning of a token if  $c \in \bigcup \alpha_i$ , where  $\alpha_i$  is the set of all non-empty initial segments of elements of  $L_i$ :

$$\alpha_i := \{u \mid \exists v. u \cdot v \in L_i, u \neq \emptyset\}$$

Given  $c \in \bigcup \alpha_i$ , function  $\Phi$  is applied to recognize the token beginning with  $c$ .

$$\begin{aligned} \Phi h \cdot s &= \text{if } s = \langle \rangle \text{ then } \langle \rangle \text{ else} \\ &\quad \text{if } \text{hd } s \in \bigcup \alpha_i \text{, then } (\Phi h \cdot (\text{hd } s)) \cdot (\text{tl } s) \#^1 \text{ else} \\ &\quad \Psi h \cdot (\text{tl } s) \end{aligned}$$

$\Phi \in H \rightarrow Str \rightarrow Str \rightarrow (Tk^* \times H)$  applied to two strings  $s_1$  and  $s_2$  finds the longest initial substring of  $s_1 \cdot s_2$  which is in some  $\alpha_i$ .  $\Phi$  repeatedly appends the first character of  $s_2$  to  $s_1$  until  $s_1 \in \bigcup \alpha_i$  ( $s_1$  is the beginning of a token) and  $s_1 \cdot (\text{hd } s_2) \notin \bigcup \alpha_i$  (i.e.  $s_1$  is the longest possible substring). In this situation  $\Phi$  returns the token corresponding to  $s_1$  concatenated with the tokens given by the remainder ( $s_2$ ).

$$\begin{aligned} \Phi h \cdot s_1 \cdot s_2 &::= \text{if } s_2 = \langle \rangle \text{ then } \langle \rangle \text{ else } (S_i \cdot h \cdot s_1) \#^1 \text{ else } \text{error} \text{ else} \\ &\quad \text{if } s_1 \cdot (\text{hd } s_2) \in \bigcup \alpha_i \text{, then } \Phi h \cdot (s_1 \cdot (\text{hd } s_2)) \cdot (\text{tl } s_2) \text{ else} \\ &\quad \text{if } s_1 \in L_i \text{, then } (S_i \cdot h \cdot s_1) \#^1 \cdot (\Psi (S_i \cdot h \cdot s_1) \#^2 \cdot s_2) \text{ else } \text{error} \end{aligned}$$

### 3.2. Micro syntax of $LS$

To define the micro syntax we have to specify the domains  $Tm$ ,  $Vl$ , languages  $\{L\}$ , semantic functions  $\{S\}$ , and the initial identifier table  $s_0$ . For  $LS$  we have 5 regular languages.

$L_{id}$  is the set of all identifiers. The corresponding semantic function  $S_{id}$  will detect reserved words and encode all other identifiers as unique integers, using the identifier table  $h$ .

$L_n$  is the set of numbers.  $S_n$  defines which integer value is denoted by a given string of digits.

$L_d$  is the set of delimiters consisting of just one character.  $S_d$  will map each character into its representation as a token.

$L_p$  is the set of all tokens beginning with a period, i.e.  $\{".", "..."\}$ .

$L_c$  is the set of all tokens beginning with a colon, i.e.  $\{":", "::="\}$ .

The complete formal definition of the scanner for  $LS$  can be found in appendix 1.

### 4. Syntax

The syntax of  $LS$  is defined by a context free grammar. We extend the definition of grammars to include unique labels for each production. We require that the grammar be unambiguous. In this case the grammar defines a mapping from token sequences  $Tk^*$  to parse trees. The production labels become labels of parse trees.

#### 4.1. Labeled context free grammars

Ordinarily a grammar is given by a set of productions. Instead, we specify labeled grammars by a set of tables  $L$  together with a mapping  $P$  from labels to productions.

$n \in Nt$	Nonterminal symbols
$Tm$	Terminal symbols
$\lambda \in L$	Labels
$u, v, w \in W = (Nt \cup Tm)^*$	Words
$P \in L \rightarrow Nt \times W$	Productions

$s_0 \in Nt$  is a distinguished start symbol. If  $P\lambda = (n, w)$  we write

$$[\lambda] n::=w.$$

A grammar  $G$  is completely specified by the elements  $Nt$ ,  $Tm$ ,  $L$ ,  $P$ , and  $s_0$ .

#### 4.1.1. The accepted language

On  $W$  we define a binary (infix) relation  $\rightarrow \in W \times W$  as follows.

$$u \cdot (n) \cdot v \rightarrow u \cdot w \cdot v \text{ iff } \exists \lambda \in L. P(\lambda) = (n, w).$$

Let  $\rightarrow^*$  be the reflexive, transitive closure of  $\rightarrow$ . For a grammar  $G$  we define

$$L(G) = \{w \in Tm^* \mid s_0 \rightarrow^* w\}.$$

Note, that  $L(G)$  is a language over strings of terminal symbols.

Recall that the domain of tokens was defined as  $Tk = Tm \times Vl$ . We extend  $L(G)$  for strings of tokens and define:  
The language accepted by a labeled context free grammar  $L_V(G)$  is

$$L_V(G) = L(G) \times Vl^*$$

Here  $\times^*$  is the obvious componentwise extension of the cartesian product to sequences:

$$S \times^* T = \{(s_1, t_1), \dots, (s_n, t_n) \mid s_1, \dots, s_n \in S, t_1, \dots, t_n \in T\}$$

#### 4.1.2. Parse trees

A tree  $s \in S = W_L(Tk)$  is a partial parse tree if

- whenever  $\hat{s} = (\lambda \tau_1 \dots \tau_m)$  is a subtree of  $s$ , then  $P(\lambda) = (n, z_1 \dots z_m)$  and for all  $0 \leq i \leq m$  we have either  $\tau_i \in Tk$  and  $\tau_i \#_1 = z_i$  or  $\text{root } \tau_i = z_i$ .

A partial parse tree  $s$  is a parse tree if in addition

- $\text{root } s = (s_0, w)$  for some  $w$ .

#### 4.1.3. The function defined by a labeled grammar

We require that a labeled context free grammar be unambiguous; that is, whenever  $s_1, s_2 \in W_L(Tk)$  then  $\text{leaves } s_1 = \text{leaves } s_2$  implies  $s_1 = s_2$ .

## Syntax

67

Each unambiguous labeled context free grammar  $G$  defines a “parsing function”  $\text{parse} \in Tk^* \rightarrow St$  as follows.

$$\text{parse}(w) = \text{if } w \in L(G) \text{ then } s, \text{ such that } \text{leaves}(s) = w \text{ else error}$$

where  $s$  is a parse tree for  $G$ .

### 4.2. Syntax of $LS$

Using the above theory, the syntax of  $LS$  is defined by giving a labeled context free grammar for  $LS$ . The function  $\text{parse}$  for  $LS$  is the function defined by this grammar. A labeled context free grammar for  $LS$  is given in appendix I. This grammar also satisfies the SLR(1) condition [De71]; we will use this property in the construction of a parser for  $LS$  which is described in chapter IV.

## 5. Tree transformations

Tree transformations define a mapping from parse trees into abstract syntax. We first define the abstract syntax of  $LS$  and then specify tree transformations for  $LS$ .

### 5.1. Abstract syntax

The abstract syntax  $asyn$  is the sum of a set of syntactic domains. These domains define syntactic entities of  $LS$ , such as identifiers, expressions, statements and so on. We use a slightly different notation here than we use for semantic domains. For example instead of writing

$$\Theta \in Stm = (Exp \times Exp) + (Exp \times Com) + \dots$$

we write

$$\Theta ::= E_0 := E_1 \mid \text{while } E \text{ do } od \mid \dots$$

for  $E_i \in Exp$  and  $\Gamma \in Com$ . Thus, this notation not only defines the domains but also a particular representation which is used to refer to elements of these domains.

#### 5.1.1. Syntactic domains

$Pgm = Stm$  Programs  
 $\Omega \in Op$  dyadic operators (not further defined here)

## III. Source and target languages

68

$O \in Mop$  monadic operators (not further defined here)  
 $I \in Id$  Identifiers (not further defined here)  
 $N \in Num$  Numerals (not further defined here)  
 $B \in Pgm$  Programs  
 $E \in Exp$  Expressions  
 $\Theta \in Stm$  Statements  
 $\Gamma \in Com$  Commands  
 $T \in Typ$  Types  
 $\Delta \in Decl$  Declarations  
 $\Delta_c \in Cdef$  Constant Definition  
 $\Delta_t \in Tdef$  Type definition  
 $\Delta_v \in Vdec$  Variable declaration  
 $\Pi \in Par$  Parameters

Expressions  
 $E ::= I \mid O \ E \mid E_0 \ \Omega \ E_1 \mid E \uparrow \mid I(E^*) \mid E_0[E_1] \mid N \mid E \cdot I$

Commands  
 $\Gamma ::= N:\Theta \mid \Theta \mid \Gamma_0; \Gamma_1$

Statements  
 $\Theta ::= E_0 := E_1 \mid \text{if } E \text{ then } \Gamma_0 \text{ else } \Gamma_1 \mid f_i \mid \text{dummy} \mid$   
 $\text{while } E \text{ do } \Gamma \text{ od} \mid \text{repeat } \Gamma \text{ until } E \mid \text{goto } N \mid I(E^*) \mid$   
 $\Delta^* \begin{cases} \text{begin } \Gamma \text{ end} \\ \text{begin } \Gamma \text{ end} \end{cases}$

Declarations  
 $\Delta ::= \text{const } \Delta^* \mid \text{type } \Delta^* \mid \text{var } \Delta^*_v \mid$   
 $\text{procedure } I(\Pi^*); \Theta \mid \text{function } I(\Pi^*); T; \Theta$

Types  
 $T ::= I \mid (I_1, \dots, I_n) \mid E_1 \cdot E_2 \mid \text{array}[T_1] \text{ of } T_2 \mid$   
 $\text{record } I_1; T_1; I_2; T_2; \dots; I_n; T_n \text{ end} \mid \uparrow I$

Constant Definitions  
 $\Delta_c ::= I = E$   
 Type Definitions  
 $\Delta_t ::= I = T$   
 Variable Declarations  
 $\Delta_v ::= I:T$

**Parameters**

$$\Pi ::= I_1; I_2 \mid \text{var } I_1; I_2$$

In addition to this mathematical notation a machine readable representation is necessary for mechanical procs. This is provided by constructor and selector functions which are defined in appendix 1.

The domains  $I_d$  and  $N_{\text{um}}$  are left undefined. The structure of identifiers is part of the definition of the micro syntax. All that is relevant for the semantics of  $LS$  is that we can test identifiers for equality. Similarly, the structure of numerals is irrelevant. We are only interested in the value denoted by a numeral.

**5.2. Tree transformations for  $LS$** 

A tree transformations is defined mapping parse trees into elements of the abstract syntax. The function  $\mathcal{E}$  is defined recursively on parse trees. One definitional clause is provided for possible node (label of the grammar). For example, the grammar of  $LS$  contains the following productions:

$$\begin{aligned} [STMT\_5] \quad STMT &::= \text{ifsymbol EXPRESSION symbol COM ifsymbol} \\ [STMT\_6] \quad STMT &::= \text{ifsymbol EXPRESSION symbol COM} \\ &\quad \text{elsesymbol COM ifsymbol} \end{aligned}$$

$$[FACT\_1] \quad FACT ::= (\text{parenthesis symbol EXPRESSION parenthesis symbol})$$

The corresponding clauses of the tree transformation are

$$\begin{aligned} \mathcal{E}(STMT\_5, r_1, r_2, r_3, r_4, r_5) &= mkeif / \mathcal{E}(r_2, \mathcal{E}(r_4), mkeastmt(mkendum)) \\ \mathcal{E}(STMT\_6, r_1, r_2, r_3, r_4, r_5, r_6, r_7) &= mkeif / \mathcal{E}(r_2, \mathcal{E}(r_4, \mathcal{E}(r_6))) \\ \mathcal{E}(FACT\_1, r_1, r_2, r_3) &= \mathcal{E}(r_2) \end{aligned}$$

For example, in the abstract syntax there is no **if** statement without **else** part. The first clause above maps a conditional without matching **else** part into an **if** statement with empty **else** part. The third line shows that parenthesis in expressions are ignored in the abstract syntax. The complete definition of  $\mathcal{E}$  can be found in appendix 1.

**6. Semantics of  $LS$** 

The semantic definition of  $LS$  follows [Te77a]. In particular we use Tennents separation of dynamic and static semantics. The intuition is, that a function is static if it can be evaluated at compile time, dynamic otherwise;

but formally this division is arbitrary, since there is no notions of compile and run time in a formal semantics.

The dynamic semantics of  $LS$  is specified on a fairly low level using concepts well known to the compiler constructor. Clearly, a more "abstract" definition of  $LS$  is possible. In this case the equivalence with a lower level description such as ours has to be proven as part of the compiler verification. Proof techniques to do this are well established (see for example [MS76]).

Since we are more interested in techniques of specification and verification rather than semantics theories we choose to base the compiler proof on a more detailed definition. This definition can be used directly as the formal basis for the compiler.

The formal definition of  $LS$  is presented in appendix 1. In this section we explain some of the interesting points of this definition. We will be informal in our explanations and relate formal objects of the definition to concepts known to compiler constructors.

**6.1. Semantic concepts****6.1.1. Semantic domains**

In this section we describe the semantic domains used in the definition of the semantics of  $LS$ .

$T$  and  $N$  are truth values and integers respectively.  $T$  is used to model predicates of the descriptive language; the data type boolean is not related to  $T$ .

$I = \{int_{-\infty}, \dots, int_{+\infty}\} \subseteq N$  is the set of index values. These are those integers that can be represented in our target machine. The specific values of  $int_{+\infty}$  and  $int_{-\infty}$  are irrelevant; we assume however, that  $0, 1 \in I$ .  
 $Tg = \{int, bool, \nu_1, \dots, \nu_i, \dots\}$  is a flat domain of type tags. Type tags are used to create new names for types in  $LS$ , i.e. as in Pascal each new type definition creates a unique type. Note, that we cannot simply use the identifier provided in the  $LS$  program to name types since there are types without explicit name and scope rules would cause ambiguities.

The concept of locations in  $LS$  is fairly complex. First we have a set of absolute locations. These can be thought of as addresses of a real physical memory. Abstract locations are divided into two distinct sets  $L_d$  and  $L_s$ , called dynamic and static locations respectively. The difference between  $L_d$  and  $L_s$  is that dynamic objects (created with the **new** statement) are assigned

dynamic locations (on the heap). Variables declared in a block or procedure are allocated static locations (on the stack).

A location in  $L_d$  or  $L_s$  is able to hold a single object (such as an integer or boolean value). To be able to describe addresses assigned to complex objects such as record and arrays a more complex structure of locations is required. We call those more general locations  $L$ -values ( $L_v$ ). Again, we distinguish static and dynamic  $L$ -values,  $L_{v_s}$  and  $L_{v_d}$ . We define

$$\begin{aligned} v \in L_{v_s} &= L_s + (Id \rightarrow L_{v_s}) \vdash (I \rightarrow L_{v_s}) \\ v \in L_{v_d} &= L_d + (Id \rightarrow L_{v_d}) + (I \rightarrow L_{v_d}) + \{NIL\} \end{aligned}$$

This means that a static  $L$ -value is either a single location ( $L_s$ ) or the address of a record ( $Id \rightarrow L_{v_s}$ ) or the address of an array ( $I \rightarrow L_{v_s}$ ). The address of a record is a mapping which given a component identifier returns the address of this component. Similarly, the address of an array is a mapping from index values into addresses of the components of the array.

In addition to static and dynamic locations we introduce relative locations  $L_r$  and relative  $L$ -values  $L_{v_r}$ . Since procedures and functions in LS can be recursive different incarnations of the same procedure (or function) give raise to different absolute locations for the local variables of this procedure. To describe the meaning of a procedure statically, i.e. independently of a particular call, we assign relative locations and  $L$ -values to local variables. At run time these relative locations are mapped into absolute locations by a memory frame. We define memory frames as

$$f \in F = L_r \rightarrow L_s.$$

Note, that a relative location is always mapped into a static location since there are no variables that are assigned dynamic locations.

Although memory frames are defined for  $L_r$  only they immediately extend to relative  $L$ -values as follows. ( $\hat{f}$  is the extended function)

$$\begin{aligned} \hat{f}a &= \text{if } a \in L_r \text{ then } f_a \text{ else} \\ &\quad \text{if } a \in (I \rightarrow L_{v_r}) \text{ then } \lambda I. \hat{f}(a) \text{ else} \\ &\quad \text{if } a \in (Id \rightarrow L_{v_r}) \text{ then } \lambda Id. \hat{f}(a[Id]) \end{aligned}$$

We assume that there are infinitely many frames  $f_i \in F$  and that  $f_i$  are totally ordered. The function succ gives the successor of a frame, i.e.  $\text{succ } f_i = f_{i+1}$ . Further, the ranges of two different frames are distinct, i.e. for all  $f_i, f_j$  and  $a_1$  and  $a_2$  we have that  $f_i \neq f_j \Rightarrow f_i.a_1 \neq f_j.a_2$ .

Values  $V$  are objects that can be stored in memory and can be result of an expression. The only values in LS are index values and pointers, i.e. we have

$V = I + L_v$ . In particular there are no array or record values; consequently, arrays and record as whole structures cannot be assigned or passed as value parameters.

A store can be thought of as the memory of an abstract machine; it is a mapping from locations to values. To characterize a computation state completely, a store in LS has two components representing the current input and output file. We define

$$S = (L_s + L_d) \rightarrow (V + \{\text{unused}\}) \times N^* \times N^*$$

A dynamic location which is unused is mapped into unused. The use of static locations will be recorded in the environment (see below).

An answer of a LS program can either be an output file  $\in N^*$  or a runtime error, thus we have  $A = N^* + E_R$ .

Continuations  $C = S \rightarrow A$  will be referred to as dynamic continuations. Instead of introducing expression continuations we define dynamic generalized continuations as follows:

$$G_d = C + (V \rightarrow G_d)$$

It is easy to see that  $G_d$  is isomorphic to

$$V \rightarrow \dots \rightarrow V \rightarrow C$$

for any number of  $V$ 's; alternatively we can write  $V^* \rightarrow C$ . The use of dynamic generalized continuations will greatly simplify the definition of meaningful functions. In addition generalized continuations have an obvious realization in a compiler. Note, that  $V^* \rightarrow C$  can be interpreted as an expression continuation which takes a "stack" ( $\in V^*$ ) as argument. As can be seen later when we give the definition of meaning functions this stack is exactly the stack used to evaluate expressions in a compiler.

Each call to a procedure requires that new memory be allocated; an executing procedure together with its local memory is called an incarnation of this procedure. Addresses used in the code of a procedure are relative addresses ( $\in L_r$ ). To access the content of a relative locations it has to be converted to an absolute locations using a frame.

The incarnation of the procedure that called the currently executing procedure  $p$  is said to be the dynamic predecessor of the current incarnation of  $p$ . Every procedure has an associated lexical level, the main program has lexical level 0. If a procedure  $p$  with lexical level  $n$  is called, then at least one incarnation of all procedures with smaller lexical level ( $< n$ ) which contain

the declaration of  $p$  has to exist. If  $q$  is such a procedure with lexical level  $m$ ,  $m < n$ , then the most recent incarnation of  $q$  is said to be the static predecessor of  $p$  of level  $m$ .

If a variable with relative location  $\alpha$  is defined on lexical level  $m$  and is accessed in a procedure  $p$  with lexical level  $n$ ,  $m < n$ , then the frame belonging to the static predecessor of  $p$  of level  $m$  has to be used to convert  $\alpha$  into an absolute location. We introduce a domain of displays to describe this use of frames. The concept of a display is well known to the compiler constructor, see for example [RR64, Gr71].

We define

$$X \in X = (B \rightarrow X) \times (B \rightarrow F) \times X \times B.$$

Here  $B = N$  is the domain of lexical levels. A display  $X$  has the following components

- A mapping from lexical levels to displays which gives the static predecessor of the current display.

- A mapping from lexical levels to frames giving the frame corresponding to the static predecessors.

- The dynamic predecessor.

- The lexical level of the currently executed procedure.

At any point the meaning of the remainder of the program depends on the current display. Therefore, we introduce generalized continuations  $G$  which depend on a particular display:

$$G = X \rightarrow G_d.$$

The meaning of a procedure is a mapping from a list of parameter values and a continuation into the domain of continuations. We could define  $V^* \rightarrow G \rightarrow G$ . But this domain is isomorphic to  $G \rightarrow V^* \rightarrow G$  which in turn is isomorphic to  $G \rightarrow G$ . Thus we define  $P = G \rightarrow G$  as the domain of procedure values. Functions can be modeled by the very same domain. We do not distinguish procedure and function values.

We use two kinds of environments in the definition of L.S. Static environments  $U_s$  contain information about identifiers and numbers relevant for static semantic analysis. Environments  $U$  are used to describe dynamic semantics.

In

$$U_s = (Id \rightarrow M_d) \times (Num \rightarrow T) \times (Tg \rightarrow T)$$

each identifier is mapped into a mode; each number is mapped to a truthvalue indicating whether or not this number is a defined label. Type tags  $Tg$  are

mapped into truthvalues to indicate which of the tags have been used to define a type; this is necessary to be able to generate a virgin type tag for the next type declaration.

An environment maps each variable identifier into a relative L-value, each function and procedure identifier into a procedure value and each label (numerical) into a continuation. In addition for each identifier and label the lexical level of its definition is given. Finally, a component  $L_r \rightarrow T$  keeps track of used relative locations.

$$\begin{aligned} \rho \in U = & (Id \dots Lv_r) \times \\ & (Id \rightarrow P) \times \\ & (Id \rightarrow \delta) \times \\ & (Num \rightarrow G) \times \\ & (Num \rightarrow B) \times \\ & (L_r \rightarrow T) \times \\ & B \end{aligned}$$

### 6.1.2. Types and modes

The domain  $Ty$  of types provides an "internal representation" for user defined types. We use a definition similar to that of the abstract syntax. A  $r \in Ty$  is enclosed in  $[ \dots ]$ , similar to elements of the abstract syntax which are enclosed in  $[ \dots ]$ .

Each  $r \in Ty$  has various fields, one of these is a type tag that uniquely identifies the type. An exception is the *nil*/type which is a special pointer type. It is special because *nil* is common to all pointer types and has to be treated differently during type checking. In the following definition products are not strict! The domain  $Ty$  contains infinite objects (e.g. recursive types).

$$Ty = [\nu: sub, i_1:i_2] \mid [\nu: \dagger r] \mid [nil] \mid [\nu: array, r_1:r_2] \mid [\nu: record, \{f_1:f_1, \dots, f_n:f_n\}]$$

Each identifier and expression in L.S is assigned a mode. For example, a mode specifies, whether an object is a constant, a variable, a value, a procedure and so on.

A mode  $\mu \in M_d$  is enclosed in  $[ \dots ]$ ; we have the following modes:

- $[var:r]$  a variable of type  $r$ .
- $[val:r]$  a var parameter inside a procedure. Its value is a variable of type  $r$ .
- $[val:r]$  a value of type  $r$ .
- $[type:r]$  a type identifier denoting type  $r$ .

- $[const.r]$  a constant of type  $r$  with value  $r$ .
- $[proc.\mu_1, \dots, \mu_n]$  a user defined procedure with parameter modes  $\mu_1, \dots, \mu_n$ .
- $[sproc]$  a special procedure. A special procedure is a predefined procedure. These are treated differently from user defined procedures to allow for more flexible parameter checking (e.g. predefined procedures can be generic, new is an example).
- $[afun.\mu_1, \dots, \mu_n;r]$  an active function with parameter modes  $\mu_1, \dots, \mu_n$  and result type  $r$ . A function is active inside its own body; it is passive otherwise. The distinction is relevant since it is possible to assign to an active function but not to a passive function.
- $[pfun.\mu_1, \dots, \mu_n;r]$  a passive function.
- $[sfun]$  a special function (see remark under special procedure).

### 6.1.3. Auxiliary functions, static semantics

The following is a list of functions used in defining the statics semantics of LS with a brief description.

$distinct \in D^* \rightarrow T$  is defined for any flat domain  $D$  for which an equality test is defined. It is used to test lists of numbers and identifiers for distinctness.

$typetag \in Ty \rightarrow Ty$  selects the type tag of a type.

A type is an index type if it is a subrange (note, that integers are a subrange  $int_{-\infty}^{+\infty}$ ).  $isindex \in Ty \rightarrow T$  returns  $TT$  for index types.

A type is returnable if it can be the result type of a function, we have  $isreturnable \in Ty \rightarrow T$ . In LS we allow index types, pointers, and the nil type to be returned by functions

$overlap, contains \in Ty \rightarrow Ty \rightarrow T$  for two subranges these functions determine whether they overlap or whether one is contained in the other. They are useful in the definition of assignments: two subranges that do not overlap cannot be assigned; if the right hand side is contained in the left hand side, no runtime check is required.

$union \in Ty \rightarrow Ty \rightarrow Ty$  constructs the smallest subrange that contains two subrange types.

$type \in Md \rightarrow Ty \rightarrow Ty$  selects the type of a mode. For example type  $[var.r] =$

$r$ , type is undefined for procedure modes and special functions.

$equal \in Ty \rightarrow Ty \rightarrow T$  tests two types for equality.

$compatible \in Ty \rightarrow Ty \rightarrow T$ , two types are compatible if there is a possible runtime check that allows these types to be assigned to each other.

- For example, equal types are compatible, two subranges that overlap are compatible, a pointer type is compatible with the nil type.
- $Assignable \in Md \rightarrow Md \rightarrow T$ ; mode  $\mu_2$  is assignable to mode  $\mu_1$  if  $\mu_1$  is a variable, if types of  $\mu_1$  and  $\mu_2$  are compatible, and if the type of  $\mu_2$  is returnable. LS only allows assignment of objects that are also returnable as function results. Arrays or records cannot be assigned.
- If an expression has mode  $\mu_2$ , passable  $\in Md \rightarrow Md \rightarrow T$  determines if it can be passed to a formal parameter that has mode  $\mu_1$ . For var parameters LS requires type equality, in particular, subranges have to match exactly. For value parameters the two modes have to be assignable, i.e. complex data objects cannot be passed as value parameters.
- $isvar, isval, \in Md \rightarrow T$  check whether a mode is a value or a variable.
- $isbool, isint \in Ty \rightarrow T$  test if a type is integer or boolean.
- $sp \in Id \rightarrow Md^* \rightarrow T$  checks a call to a special procedure for validity.

Given the name of the special procedure and the modes of the actual parameters  $sp$  returns  $TT$  if this call is valid. Note, that this is a very flexible schema, for example a procedure  $print$  could be made generic for all types; also print could take a variable number of parameters for each call. We assume, however, that  $print$  has only one integer argument; extensions are easy to add to our compiler.

$sf \in Id \rightarrow Md^* \rightarrow Md$  is similar to  $sp$  but checks special functions instead. The result of  $sf$  is the mode of the object returned by the special function.

$true, mode \in Md$  are modes of the constants  $true$  and  $false$ .

$inconst \in I \rightarrow Md$  constructs an integer constant of a given value.

$integer \in Ty$  is the type integer.

$boolean \in Ty$  is the type boolean.

$memvar, mode \in Md \rightarrow Md$  makes a value mode with the same type as a given mode.

$w \in \Omega \rightarrow Md \rightarrow Md \rightarrow Md \rightarrow O \rightarrow O$ ,  $w \in \Omega$  is the mode of the object returned by the special function.

### 6.2. Static semantics

#### 6.2.1. Declarations

A list of declarations is checked by  $d$  for its correctness.  $d$  also updates the environment to reflect the new declarations. We have  $d \in Dec \rightarrow U_s \dashv U_s$ . To

treat constant, type, and var declarations  $d$  uses the functions  $dc \in Cdef \rightarrow U_s \rightarrow U_s$ ,  $dt \in Tdef \rightarrow U_s \rightarrow U_s \rightarrow U_s$ , and  $dv \in Vdec \rightarrow U_s \rightarrow U_s$ .

One of the more interesting points in  $d$  is the treatment of recursive types. In  $d$  we have the clause

$$d[\![type\Delta_1; \dots; \Delta_n]\!]_S = fix(dt[\![\Delta_1; \dots; \Delta_n]\!]_S).$$

Note, that  $dt$  and  $dt^*$ <sup>4</sup> require two environments as argument. To get a feel for the meaning of the fixed point, consider the limit case, where we have  $d[\![type\Delta_1; \dots; \Delta_n]\!]_S = x$  such that  $x = dt[\![\Delta_1; \dots; \Delta_n]\!]_S$ . If any of the  $\Delta_i$ 's contains a pointer type, the reference will be resolved by looking into  $x$  instead of  $S$ , whereas  $S$  is used for all other identifiers: they have to be defined in a global scope.

Another interesting case is that of procedure definitions.

$$d[\![procedure I(\Pi_1, \dots, \Pi_n); \Theta]\!]_S = let ((\mu_1, \dots, \mu_n), \varsigma_n) : p [\![\Pi_1, \dots, \Pi_n]\!]_S in$$

$$\quad if distinct([\![\Pi_1, \dots, \Pi_n]\!], [\!\Theta]\!) then$$

$$\quad \quad if q[\!\Theta\!]_{S_n} [\![proc(\mu_1, \dots, \mu_n)]\!]_I then$$

$$\quad \quad \quad S_n [\![proc(\mu_1, \dots, \mu_n)]\!]_I$$

$p$  is used to determine the modes of the parameters as well as the environment  $S_n$  resulting from  $S$  after entering the parameters with their modes. If all identifiers used as parameters and in the procedure body are distinct and if the procedure body is semantically correct (see definition of  $q$ ) in  $S_n$ , then the resulting environment is obtained by entering the corresponding procedure mode for  $I$  in  $S_n$ .

The remaining cases are straightforward.

### 6.2.2. Types

$t \in Typ \rightarrow U_s \rightarrow U_s \rightarrow (Ty \times U_s)$  determines the meaning of a type in the program. A given type denotation is evaluated and a representation ( $\in Ty$ ) is constructed. Two points are to be observed here.  $t$  changes the environment in two ways. First, an enumeration defines a set of identifiers as constants of this enumeration types; these have to be entered by  $t$  in the environment. Furthermore, to represent new types,  $t$  uses up type tags which in turn have to be marked as being in use in the new environment.

The second point to observe is that like  $dt$ ,  $t$  too takes two environments as argument. The reason is exactly the same here:  $t$  uses the first environment in all cases except to determine the type pointed to by a pointer type declaration.

4.) see appendix 1.

### 6.2.3. Labels and identifiers

$j \in Com \rightarrow Num^*$  returns the list of all labels defined on the top level of a given command.

$k$  returns the list of all procedures and functions defined in a block.

Function  $t$  returns all identifiers defined in a given language construct. It is defined for several syntactic domains:  $i \in (Par + Stmt + Dec + Typ) \rightarrow Id^*$ .

### 6.2.4. Expressions

The function  $e \in Expr \rightarrow U_s \rightarrow Md$  evaluates an expression in a given static environment and returns the mode of the expression. In computing this mode  $e$  checks if all declaration and type constraints are satisfied. If this is not the case  $e$  will return  $\perp$ . We will explain some typical clauses of the definition.

For an integer constant we have  $e[\![N]\!]_S = [const.n.[int;abs;n][n]]$ , where  $n$  is the value of the constant  $N$  (given as the token value). That is, the mode of a constant is the constant mode with the value being the value of the constant and the type being the smallest possible subrange containing this constant.

The mode of a binary operator is given by the function  $w$  applied to the modes of the operands:  $e[\![E_0 \Omega E_1]\!]_S = w[\Omega](e[\![E_0]\!]_S, e[\![E_1]\!]_S)$

For an indexing operation to be well typed  $e$  has to check that

- the indexed expression is of an array type,
- the type of the indexing expression is compatible with the index type of the array type.

The resulting mode is a var in any case: there are no values or constants of array type. The corresponding definitional clause is:

$$e[\![E_0[E_1]]\!]_S = \text{if } [var.array\ r_1\ r_2].type\ e[\![E_0]\!]_S \text{ then}$$

$$\quad \quad \quad \text{if compatible}\ [type\ e[\![E_1]\!]_S] \text{ then}$$

$$\quad \quad \quad \quad \quad \text{if isvar } e[\![E_0]\!]_S \text{ then } [var;r_2]$$

$$\quad \quad \quad \quad \quad \text{if isup } e[\![E_0]\!]_S \text{ then } [var;r_2]$$

Other expressions are defined in a straightforward way following the above schema.

### 6.2.5. Statements

Statements are checked by  $g \in Stmt \rightarrow U_s \rightarrow T$  for their semantic validity. The definition is straightforward, for example assignments are defined as

$$g[\![E_1 := E_2]\!]_S = assignable[e[\![E_1]\!]_S](e[\![E_2]\!]_S)$$

That is, an assignment is legal if the modes of left and right hand side are assignable.

A slightly more complex example are while loops. There are two semantic restrictions: the test must yield a boolean result and the body of the while statement must be valid. For the latter test  $c$  is used to evaluate a command. Note, that  $c$  is passed a modified environment in which all labels in the while body are defined; this effectively makes the body of the while loop a new scope for the purpose of labels. It becomes impossible to jump into a while body.

$$\begin{aligned} q[\text{while } E \text{ do } \Gamma \text{ od}]_c &= \text{isbool}(\text{type } c[E]) \wedge \\ &\quad \text{distinct}(c[\Gamma]) \wedge \\ &\quad c[\Gamma]_s[\{\dots, \text{true}, \dots\}/j[\Gamma]] \end{aligned}$$

### 6.2.6. Commands

Commands are lists of labeled or unlabeled statements. A list of statements is semantically valid, if all individual statements are valid. Thus for  $c \in \text{Com} \rightarrow U_s \rightarrow T$  we have the simple definition:

$$\begin{aligned} c[\text{N}:\Theta]_s &= g[\Theta]_s \\ c[\Theta_1:\Theta_2]_s &= c[\Theta_1]_s \wedge c[\Theta_2]_s \\ c[\Theta]_s &= g[\Theta]_s \end{aligned}$$

### 6.3. Dynamic semantics

Meaning functions of the dynamic semantics are denoted by capital script letters  $\mathcal{E}$ ,  $\mathcal{D}$  and so on.

#### 6.3.1. Auxiliary functions

By introducing generalized continuations operations on continuations  $C$  have to be extended to operate on  $G_d$ . Suppose we are given  $\gamma \in G_d = V \rightarrow \dots \rightarrow V \rightarrow S \rightarrow A$ . We want to define a continuation  $\dot{\gamma}$  which reflects a change  $f \in S \rightarrow \gamma$  to the store. I.e. we want

$$\gamma e_1 \dots e_n(f \sigma) = \dot{\gamma} e_1 \dots e_n \sigma$$

The function  $\text{change} \in G_d \rightarrow (S \rightarrow S) \rightarrow G_d$  computes  $\dot{\gamma}$  given  $f$  and  $\gamma$ .

$\text{update} \in G \rightarrow G$  describes the effect on a continuation of assigning a value to a location:

$$\text{update } \gamma = \lambda X \alpha. \text{change}(\gamma X)(\lambda \sigma. \sigma[\epsilon/\alpha])$$

$\text{content} \in G \rightarrow G$  describes the effect of accessing a location in memory. Let  $\gamma$  be a continuation that takes  $n$  arguments  $\in V$ , then  $\text{content}$  is defined such that

$$\text{content } \gamma = \lambda X a_1 \dots e_n \sigma. (\gamma X)(\sigma a_1 \dots e_n \sigma)$$

$\text{cond} \in G \rightarrow G \rightarrow G$  defines a "conditional continuation". Given arguments  $\gamma_1$  and  $\gamma_2$   $\text{cond}$  defines a continuation which takes an argument  $\epsilon$  and depending on the value of  $\epsilon$  applies either  $\gamma_1$  or  $\gamma_2$ . We define  $\epsilon$  to be true if  $\epsilon \mid I$  is zero.

$$\text{cond} \gamma_1 \gamma_2 = \lambda X \epsilon. \text{If } (\epsilon \mid I) = 0 \text{ then } \gamma_1 X \text{ else } \gamma_2 X$$

$\text{binop} \in \text{Op} \rightarrow G \rightarrow G$  and  $\text{unop} \in \text{Mop} \rightarrow G \rightarrow G$  define the meaning of binary and unary operators respectively

$\text{SP} \in Id \rightarrow M^*$ ,  $I \rightarrow G \rightarrow G$  defines the meaning of special procedures. In our case we have  $\text{print}$  and  $\text{new}$ . Similarly,  $Sf \in Id \rightarrow G \rightarrow G$  defines the meaning of special functions; in LS we have  $\text{eof}$  and  $\text{read}$ .

$\text{verifys} \in J \rightarrow J \rightarrow G \rightarrow G$  Given two index values and a continuation  $\text{verifys}$  verifies that a value applied to the continuation is in the subrange defined by the two index values. If this condition is not satisfied  $\text{verifys}$  will cause a runtime error. Similarly,  $\text{verifyn} \in G \rightarrow G \rightarrow G$  will check if a value supplied to a continuation is a pointer not equal to  $\text{nil}$ .  $\text{verifyn}$  is used to ensure that no nil pointer is dereferenced.

$\text{index} \in G \rightarrow G$  and  $\text{select} \in Id \rightarrow G \rightarrow G$  describe the effect of indexing an array and selecting a component of a record.

Generalized continuation cause one slight problem in the definition of jumps. Suppose we describe the meaning of  $\text{goto } n$  where the continuation following this jump is  $\gamma$  and the continuation of the label  $n$  is  $\dot{\gamma}$ . In certain cases (e.g. if the jump leaves a function) the domains of  $\gamma$  and  $\dot{\gamma}$  may differ and require different numbers of arguments  $\in V$ . Thus, before jumping we have to *adjust* the continuation we jump to. As mentioned above, we may look the  $V$  arguments of  $\gamma$  as a stack. It is well known to the compiler constructor that jumps out of functions require adjustment of the runtime stack. This fact is reflected in the formal definition.  $\text{adjust} \in G \rightarrow N \rightarrow G$  adjusts a given continuation to take  $n$  more arguments.

$\text{args} \in G \rightarrow N$  simply determines the number of  $V$  arguments of a continuation.

#### 6.3.2. Memory allocation

We define functions to (i) allocate static memory, (ii) dynamic memory on the heap, and (iii) functions to initialize newly allocated memory. In each

case we distinguish objects that occupy one location and complex data objects stored in several locations.  $\text{new}$ ,  $\text{newl}$ , and  $\text{newu}$  allocate static locations.  $\text{new}$  allocates a single location:  $\text{new} \in U \rightarrow (Lu, \times U)$ .  $\text{newl} \in Ty \rightarrow U \rightarrow (Ly, \times U)$  takes a type as additional argument; it allocates memory for any type. Finally  $\text{newu} \in I \rightarrow I \rightarrow Ty \rightarrow U \rightarrow (Lu, \times U)$  allocates memory for arrays, given the index values which define the index range and the element type of the array.

The second set of functions is *clear* and *cleara*. These do not allocate memory, rather they initialize newly allocated memory to 0. In LS all variables are initialized to 0. Alternatives would be to leave their values arbitrary and have the program behave nondeterministically or to initialize variables to *undefined* and causing a runtime error whenever an uninitialized variable is accessed. The former requires a more complicated definition using power domains [P776, Sm78]. The latter was avoided since the resulting semantics cannot efficiently be implemented on most machines; it requires an additional bit in all memory words to indicate initialization.

The third set of functions allocates dynamic locations (on the heap). We have  $\text{heap} \in S \rightarrow (Lu_d \times S)$ ,  $\text{heapa} \in I \rightarrow Ty \rightarrow S \rightarrow (Lu_d \times S)$ , and  $\text{heapl} \in Ty \rightarrow S \rightarrow (Lu_d \times S)$  which are analogous to *new*, *newa*, and *newl* but in addition initialize the allocated memory.

### 6.3.3. Declarations

Declarations cause a change in the environment: declared variables are allocated, procedure and function values are entered in the environment. In addition declarations require the execution of initialization code.

Allocation of memory for variables is defined by the functions

$$\begin{aligned} V &\in Decl \rightarrow U_s \rightarrow U \rightarrow U \\ V^* &\in Decl^* \rightarrow U_s \rightarrow U \rightarrow U \\ V_v &\in V\text{Decl} \rightarrow U_s \rightarrow U \rightarrow U \\ V_v^* &\in V\text{Decl}^* \rightarrow U_s \rightarrow U \rightarrow U \\ D &\in Decl \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G \\ D^* &\in Decl^* \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G \\ D_v &\in V\text{Decl} \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G \\ D_v^* &\in V\text{Decl}^* \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G \end{aligned}$$

Initialization of variable is described by

Inside a procedure or function parameters are treated very much like declarations: we allocate memory for each parameter; we initialize each parameter. The difference to declarations is that one location suffices for each parameter;

furthermore, parameters are initialized to the actual parameters rather than 0. Allocation is done by  $Q \in Par \rightarrow U_s \rightarrow U \rightarrow U$  and  $Q^* \in Par^* \rightarrow U_s \rightarrow U \rightarrow U$ . Parameters are initialized by  $P \in Par \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G$  and  $P^* \in Par^* \rightarrow U_s \rightarrow U \rightarrow E^* \rightarrow G \rightarrow G$ .

$\mathcal{F} \in Decl \rightarrow U_s \rightarrow U \rightarrow G^*$  and  $\mathcal{F}^* \in Decl^* \rightarrow U_s \rightarrow U \rightarrow G^*$  return the list of procedure and function values defined by a list of declarations. The way these function and procedure values are entered in the environment is described under statements where we discuss the meaning of blocks.

Let us explain the meaning of procedure; functions are treated similarly.

We have the definition:

$$\begin{aligned} \mathcal{F} [[\text{procedure } I[\Pi_1, \dots, \Pi_n], \Theta]] \rho &= (\pi) \\ \text{where } \pi &= \lambda \gamma. \text{enter}(n+1); P[[\Pi_1, \dots, \Pi_n], \varsigma_1] \rho; \\ B [[\Theta]] \varsigma_2 \rho; \text{exit } \gamma \\ \text{where } \varsigma_2 &= \varsigma_1 [[\text{proc } \mu_1, \dots, \mu_n]/I] \\ \text{where } ((\mu_1, \dots, \mu_n), \varsigma_1) &= p[[\Pi_1, \dots, \Pi_n]] f \\ \text{where } p &= Q [[\Pi_1, \dots, \Pi_n], \varsigma_1(\text{next } \rho)] \\ \text{where } n &= level \rho \end{aligned}$$

$\varsigma_2$  is the static environment in which formal parameters are bound to their mode and the procedure identifier is bound to the correct procedure mode. The environment  $\rho_1$  has memory allocated for parameters.

The value of the procedure then is given by creating a new frame (*enter*), initializing all parameters, executing the body of the procedure in  $\varsigma_2$  and  $\rho_2$ , and discarding the frame for this procedure incarnation (*exit*).

### 6.3.4. Expressions

Expressions are evaluated by  $\mathcal{E}$ ,  $\mathcal{A}$ ,  $\mathcal{L}$ , and  $\mathcal{R}$ . Here,  $\mathcal{E}$  evaluates an expression without any coercion.  $\mathcal{A}$  takes a mode as additional parameter and coerces the expression to this particular mode.  $\mathcal{L}$  and  $\mathcal{R}$  essentially are special cases of  $\mathcal{A}$  and are used to evaluate expressions on the left hand side and right hand side of an assignment respectively. All definitions are straightforward.

$\mathcal{E}$  does constant evaluation, i.e. if an expression is of constant mode its value is this constant. Although not strictly necessary, we made constant evaluation part of the definition. If it were not, constant evaluation could be an optimization. However, care has to be taken so that for instance the evaluation of  $1/0$  causes an error at the point defined by the formal semantics. Our definition of LS states that this zero division is to be detected at compile time.

$\mathcal{E}$  and  $\mathcal{A}$  check for runtime errors such as invalid index and dereferencing of nil-pointers. This is signified by the clauses:

$$\begin{aligned} \mathcal{E}[E]s\rho\gamma &= R[E]\rho, \text{verify } \gamma \\ \text{and} \quad \mathcal{A}[E]s\rho\gamma &= \text{if } \mu::[\text{var}, r] \text{ then } L[E]\rho\gamma \text{ else} \\ &\quad \text{if } \mu::[\text{val}, r] \text{ then} \\ &\quad (\text{if } f::[\nu:\text{sub}, \epsilon] \text{ then} \\ &\quad \text{if contains } rf \text{ then } R[E]\rho\gamma \\ &\quad \text{else } R[E]\rho, \text{verify } \epsilon_0\epsilon_2\gamma \\ &\quad \text{else } R[E]\rho\gamma) \\ &\quad \text{where } f = \text{type}(e[E]) \end{aligned}$$

### 6.3.5. Statements

In LS we distinguish commands and statements. The reason is to be able to handle jumps properly. In particular we have to rule out that jumps lead into statements.

The meaning of commands is trivially defined as

$$C \in Com \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G$$

$$\begin{aligned} C[N:\Theta]s\rho\gamma &= B[\Theta]\rho\gamma \\ C[\Theta]s\rho\gamma &= B[\Theta]\rho\gamma \\ C[\Gamma_0:\Gamma]s\rho\gamma &= C[\Gamma_0]\rho; C[\Gamma_1]\rho\gamma \end{aligned}$$

A special function  $C$  defines the meaning of a command that appears as a component of a composite statement such as the then part of an if-statement. Such a command has to be evaluated in an environment where labels defined in this command are bound to proper continuations. These labels are not bound in an enclosing scope to prohibit jumps into composite statements:

$$Ct \in Com \rightarrow U_s \rightarrow U \rightarrow C \rightarrow C$$

$$\begin{aligned} Ct[\Gamma]s\rho\gamma &= C[\Gamma]\rho\gamma \\ \text{where } \dot{\varsigma} &= \{[\dots, \text{true}, \dots]/j[\Gamma]\}, \\ \dot{\rho} &= fix(\lambda \dot{\rho}. \rho [j[\Gamma]\dot{\varsigma}\rho / j[\Gamma]] \dot{\rho}) \end{aligned}$$

$J \in Com \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G^*$  returns a list of continuations corresponding to the labels defined in a command.

$$\begin{aligned} J[N:\Theta]s\rho\gamma &= (B[\Theta])s\rho\gamma \\ J[\Gamma_1; \Gamma_2]s\rho\gamma &= J[\Gamma_1]\rho; C[\Gamma_2]\rho\gamma \cdot J[\Gamma_2]\rho\gamma \\ J[\Theta]s\rho\gamma &= \{\} \end{aligned}$$

Semantics of statements ( $B$ ) is conventional, however, we should mention two points here. For jumps we have to adjust the continuation as mentioned before. This is done by

$$B[goto N]\rho\gamma = adjust(\rho[N])(args \gamma - args(\rho[N]))$$

The second point of interest are declarations. The meaning of a block is to

- initialize locally declared variables followed by
  - execution of the body in the correct environment.
- The static environment is determined as for static semantics. In the environment we have to allocate new storage for local variables and bind procedure and function values and label continuations. The latter step requires to construct a fixed point: the continuation at a label depends on continuations of other labels. Similarly, procedure and function values depend on other procedure and function values (for recursive calls) as well as label continuations. Our definition handles mutually recursive procedures.<sup>5</sup>

$$\begin{aligned} B[\Delta^*; begin \Gamma end]\rho\gamma &= D^*[\Delta^*\dot{\varsigma}; C[\Gamma]\dot{\rho}\gamma] \\ \text{where } \dot{\varsigma} &= d[\Delta]s[\dots, \text{true}, \dots]/j[\Gamma], \\ \dot{\rho} &= fix(\lambda \dot{\rho}. \rho [j[\Gamma]\dot{\varsigma}\rho / j[\Gamma]] \dot{\rho}) \end{aligned}$$

where  $\rho = V[\Delta^*]\dot{\rho}$

### 7. The target language LT

In this section we describe the target language LT. First we present the architecture of a hypothetical machine and give an informal description of the language. Strictly speaking it is not necessary to define the machine architecture; the language can be defined independent of a particular implementation. 5.) these are disallowed by the static semantics

Having a concrete machine to talk about will make out explanations easier to understand and transparent.

Following the informal discussion we present a formal denotational definition of *LT*. Only some interesting parts are discussed here; the complete semantics is presented in appendix 2.

### 7.1. A hypothetical machine

#### 7.1.1. Design decisions

Ease of translation was one important concern. Therefore, the language contains fairly "high level" concepts. However, it will be simple to verify correct translation of *LT* in a lower machine language, for instance for a register machine.

The underlying hypothetical machine is a stack machine. The key idea underlying stack architectures is to implement some components of the runtime system for Algol-like languages in hardware [Or73]. The standard run-time organization as well as the machine architecture of the BG700 (a stack machine) [BC71, BC72, BC73] combine conceptually different objects into one stack. This stack is used to (i) evaluate expression, (ii) allocate static memory, and (iii) store administrative information about return addresses of procedures, static and dynamic links of procedure invocations.

In our machine these different objects of the stack are clearly separated which simplifies our proofs. But at the same time it will not be very difficult to encode separate components of our machine in one stack if efficiency dictates so.

#### 7.1.2. Architecture

Our hypothetical machine (see figure 2) consists of the following components.

- An infinite memory capable of storing infinite precision integers.
- A stack used for expression evaluation (potentially infinite)
- a display mechanism (see below).
- one input and one output file.

Our machine knows two kinds of addresses: relative and absolute addresses, both integers. Accessing memory can only be done with absolute addresses. Relative addresses can be converted into absolute ones using information contained in the display (see below).

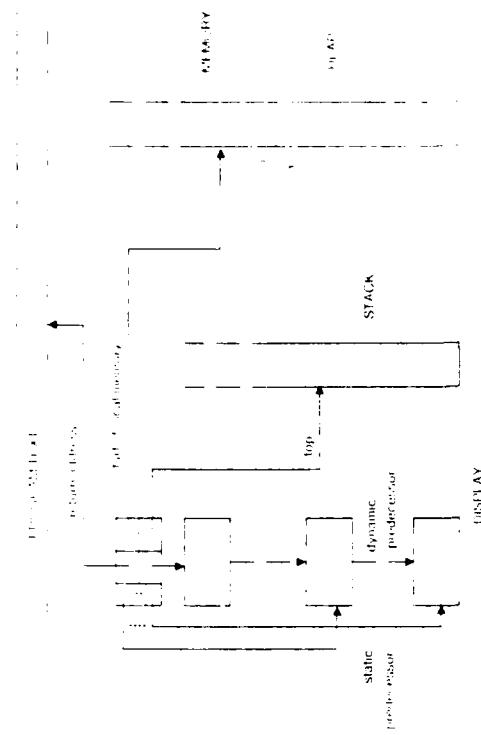


Fig. 2

The memory is divided into two parts, negative and positive addresses. The part with negative addresses will be used as heap to store dynamically allocated objects. Memory with positive addresses is administered in a stack-like fashion; i.e. memory allocated in this part becomes available when the scope in which it was allocated is left.

The display mechanism of our hypothetical machine models the domain  $X$  of the source language definition. But the domain  $D$  of target display contains more information than  $X$ . For each display we have:

- The lexical level of the currently executed procedure.
- A pointer to the displays of all static predecessors.
- The base address of each static predecessor (not shown in the picture of our machine)
- A pointer to the display of the dynamic predecessor.
- The length of the stack upon entry of the current procedure. This information is used by the *JUMP* instruction if the jump leads out of a procedure or function.
- The return address of the current procedure.

Our design of the display mechanism is not very efficient, but this is not the point of our discussion. More efficient architectures are possible that realize the same semantics.

### 7.1.3. Instructions

Let us briefly give an informal description of the instructions available on our hypothetical machine. Instructions have varying format; they may have none or several parameters. Generally, all instructions operate on the stack. That is, arguments are obtained from the stack and results are put back on the stack.

- *LIT n* pushes the constant *n* on the stack.
- *LOAD* uses the top of the stack as an address and push the content of the corresponding location.
- *ADDR n m* takes the relative address *n* and the lexical level *m* and compute aid absolute address from the current display. The result is pushed on the stack.
- *STORE* pops a value and an absolute address from the stack and updates the address with the value.
- *CALL l n m k* calls the procedure starting at label *l* defined on lexical level *n*. *m* is the number of parameters of this procedure and *k* is the number of static memory locations allocated in the calling environment. *CALL* creates a new display, stores the return address, computes the starting location of the available memory, and marks the level of the stack (current stack length minus the number of parameters).
- *EXIT* exits a procedure or function. This requires to make the dynamic predecessor of the display the new display and branching to the return address of the procedure. The stack remains unchanged, this allows a function result to be pushed on the stack before executing *EXIT*. This result is then available to the callee.
- *NCOND l* inspects the top element of the stack. If it is false ( $\neq 0$ ) then it branches to label *l*. This instruction assumes that the label is defined on the current lexical level.
- *HOP l* unconditionally branches to label *l*; *l* is assumed to be on the current lexical level.
- *JUMP l n* unconditionally branches to label *l* defined on lexical level *n*. This requires to make the display of level *n* the current display and to

adjust the stack to a length indicated in this display.

- *LABEL n* is a meta instruction and is used to define a label. Its semantics is a no-op.
- *STOP* stops program execution. This is the normal termination; "running" off the end of a list of instructions will cause an error.

- *UNOP n* performs the unary operation number *n* using the top of the stack as argument and pushing the result. We deliberately leave the precise nature of the unary and binary operations undefined. However, they have the same meaning as those in the source language. E.g.  $+^+$  in the source language has the same semantics as addition on the target machine.
- *BINOP n* performs binary operation *n* on the two top elements.

For the sake of simplicity we assume that the machine has standard input output instructions. In practice these instructions would be part of a runtime system.

- *EOF* pushes *true* (0) on the stack if the input file is empty, *false* (1) otherwise.
- *OUTP* prints the top of the stack on the output file.
- *INPT* pushes the first element of the input file on the stack.

### 7.2. Formal definition of *LT*

As mentioned earlier, no concrete syntax is required since we are not interested in parsing programs written in the target language. Also, we are not interested in checking the static semantics of *LT* programs. We merely define the dynamic semantics. If a program is not well formed, for example if a label is multiply defined, the meaning of this program will be  $\perp$ .

#### 7.2.1. Abstract syntax

The definition of the abstract syntax is straight forward. We have the domains of labels, numeral, instructions, and programs.

#### 7.2.2. Semantic domains

Integers are as usual; they model values stored in memory as well as addresses. Stacks *S* are finite lists of integers,  $S \in V$ .

Environments map labels into continuations much like in *LS*. Since there are no variables this is the only information in environments.

Displays are defined as

$$\delta \in D = N \rightarrow D \times N \rightarrow N \times D \times N \times G \times N$$

The components have the following meaning (in this order)

- Mapping from lexical levels into displays of the static predecessors.
- Mapping from lexical levels to starting addresses of the corresponding memory frame.
- The dynamic predecessor.
- The current lexical level.
- Continuation to be used upon procedure exit.
- Stack length of on entry to the current procedure.

The memory  $M = (N \rightarrow N) \times N^* \times N^* \times N$  is a mapping from addresses to values ( $N \rightarrow N$ ), the input and output file, and the first available location of the heap.

Continuation in  $L_T$  are mappings from the current machine state into answers. A machine state is characterized by the current display, the stack, and the memory. We curry and write  $G = D \rightarrow S \rightarrow M \rightarrow A$ .

The domain of answers is identical to answers of  $L_S$ ; an answer is either an output file or a runtime error.  $A = N^* + \{\text{eof}, \text{error}, \text{noatop}, \dots\}$

### 7.2.3. Semantic equations

The valuations for individual instructions are straightforward. We use functions up and down to create a new display (enter a procedure) and adjust the stack after global jumps.

Labels are handled similar as in  $L_S$ . The function  $L$  returns a list of continuations holding at labels defined in the program. These continuations are bound to the labels in the environment. As in  $L_S$  we have to construct a fixed point since the continuation at a label will depend on the continuation already bound to other labels in the environment.

### 1. Verifying a compiler

#### 1.1. The compiler

##### 1.1.1. Correctness statement

The meaning of a  $L_S$  program is given by  $smean \in Str \rightarrow M + E$ . In the previous chapter  $smean$  is formally defined as

$$smean = dsem \odot ssem \odot freestr \odot parse \odot scan.$$

Similarly, the meaning of a program in the target language  $L_T$  is given by  $tmean \in Cde \rightarrow M + E$ . The compiler to be described in this chapter will take a program in  $L_S$  and produce a sequence of code in  $L_T$ . Since there is no indeterminacy in our implementation language we can assume that the produced code depends functionally on the input program, i.e. we can look at the compiler as a function  $comp \in Str \rightarrow Cde$ .

The compiler may not terminate, in which case the result of  $comp$  is  $\perp \in Cde$ . In particular, we consider error situations as nontermination; i.e. whenever  $smean z \in E$  then  $comp z = \perp$ . We prove that whenever  $comp z$  terminates producing code  $y$ , then  $tmean y = smean z$ . Note, that this correctness statement cannot be expressed by a commutative diagram.

#### 1.1.2. Structure of the compiler

The compiler is divided into four separate programs executed sequentially. Each program has an input domain and an output domain, such that output and input of two subsequent program match. For example, the scanner produces a sequence of tokens. This sequence is input to the parser which

produces an abstract syntax tree and so on. This is the only communication between the four programs. If we were to include a sophisticated error handling mechanism more information is required, for example one might want to access the identifier table used in the scanner in later stages. Additional information of this kind can easily be communicated without invalidating our verification.

Technical details of how exactly this communication proceeds is left out of our considerations; a concrete implementation may use files or internal data structures to store intermediate forms of the program. Also, a corountining schema which avoids any intermediate storage is possible.

The division into individual modules is very conventional: we will have a scanner, a parser, a semantic analysis and a code generator. For each program we assume that *input* is the input data and *output* is the output produced. With these conventions the four components of the compiler can be formally specified as follows.

```
{true} sc {output = scan(input) | Tk^*}

{true} pa {output = treer Ø "assem(input)" | Asyn}

{true} ss {output = sem(input) | Asyn}

{true} cg {tmean(output) = dsem(input)}
```

The projection into the appropriate domains (e.g. 'Tk') enforces that the output is undefined whenever the formal definition calls for an error. Whenever an error is detected, we will assume that the corresponding program does not terminate and all subsequent programs are not executed.

The compiler is given by executing the above programs sequentially; we write informally *sc*, *pa*, *ss*, *cg*; with the understanding that each program takes the output of its predecessor as its input. Assuming the Hoare's composition rule for statements to be valid for the composition of programs, the compiler will satisfy the following specification.

```
{true}
sc; pa; ss; cg
{tmean(output) ==
dsem(ssem(treer Ø parse(scan(input) | Tk^*) | Asyn) | Asyn)}
```

Note, that the error strict functions composition degenerates to normal function composition if we project on non-error elements. Therefore, the exit

condition of the compiler is equivalent to  
 $tmean(output) = smean(input)$ .

Subsequently, we merely consider the local specification relevant for the individual programs.

### 1.2. The individual proofs

In the following sections we describe the development of the individual programs. Note, that we do not give program description or a conventional documentation of these programs. Our main interest is the development of the programs, their relation to the problem specification, and their proofs.

For each module we describe

- the theories and lemmas that had to be derived from the specifications to enable an efficient implementation,
- the basic algorithm used in the program and how its correctness can be established from the specifications and derived lemmas and theorems, and
- some relevant implementation details, representations and their effect on the correctness proof.

In addition we will focus on other interesting points of the individual program and will exemplify crucial steps in the design methodology. Some typical examples of the program text, underlying theory, and verification conditions are included in the appendix.

The programs we have to deal with are much to large to be handled by the verifier. Consequently, we have to break down the individual programs into smaller pieces that can be verified separately. However, it is necessary to include all routine, type, and variable declarations in each module to present a semantically valid program to the verifier. In total we verified 28 pieces of code comprised of 475 procedures and functions. The final program requires the proof of over 1000 verification conditions. But since each program is developed in a top-down fashion and verification takes place after each refinement step the total number of verification conditions proven during the development of the compiler is much larger.

Frequently, we use virtual data structure, types, procedures, functions, and variables. The verifier requires all these objects to be properly declared. But since virtual objects require no implementation the precise nature of their declaration is irrelevant. Thus the reader may find numerous external procedures, and meaningless type declarations such as  
 $\text{type set} = (\text{emptyset}, \text{etc})$ .

Also real objects may not be refined down to an executable level in which case these too may have nonsensical declarations.

### 2. A scanner for LS

The scanner  $\alpha_c$  has to compute the function  $\text{scan} \in Ch^* \rightarrow (Tk^* + E)$  which is defined in chapter III and appendix 1.

We prove partial correctness of the program  $\alpha_c$  with respect to the following specifications:

$\{\text{input} = \text{inputs}\} \quad \alpha_c \quad \{\text{output} = \text{scan}(\text{inputs})\}$

The definition of  $\text{scan}$  is part of the logical basis for the verification.

#### 2.1. Underlying theory

##### 2.1.1. A suitable definition

$\text{scan}$  is defined by recursive functions  $\Psi$  and  $\Phi$ . These functions specify how a set of regular languages  $L_i$  and associated semantic functions  $S_i$  are combined to define the micro syntax. For LS the languages  $L_i$  have two important properties which will allow an efficient implementation of the scanner.

- The initial segments  $\{\alpha_i\}$  of the languages  $\{L_i\}$  are identical to these languages, i.e.

$$L_{i,d} = \alpha_d = le \cdot (le + dt)^*$$

$$L_n = \alpha_n = dt \cdot dt$$

$$L_d = \alpha_d = dt$$

$$L_p = \alpha_p = pe + (pe \cdot pe)$$

$$L_e = \alpha_e = eo + (eo \cdot ==*)$$

- the  $\alpha_i$  are disjoint, i.e.  $i \neq j \rightarrow \alpha_i \cap \alpha_j = \emptyset$ .

The above properties effectively allow us to recognize one class of tokens, say identifiers, independently from other token classes. For example, once we have seen a letter in the input this can only be the beginning of an identifier, not a number for example. We make use of this property and define functions  $\Phi_i$  such that

$$\begin{aligned} \Phi \cdot h \cdot s_1 \cdot s_2 &= \text{If } s_1 \in L_i \text{ then } \Phi_{i,d} \cdot h \cdot s_1 \cdot s_2 \text{ else} \\ &\quad \text{If } s_1 \in L_n \text{ then } \Phi_n \cdot h \cdot s_1 \cdot s_2 \text{ else} \\ &\quad \text{If } s_1 \in L_d \text{ then } \Phi_d \cdot h \cdot s_1 \cdot s_2 \text{ else} \\ &\quad \text{If } s_1 \in L_p \text{ then } \Phi_p \cdot h \cdot s_1 \cdot s_2 \text{ else} \\ &\quad \text{If } s_1 \in L_e \text{ then } \Phi_e \cdot h \cdot s_1 \cdot s_2 \end{aligned}$$

where  $\Phi_i$  are defined as:

$$\begin{aligned} \Phi_{i,d} \cdot h \cdot s_1 \cdot s_2 &= \text{If } s_2 = () \text{ then } ((S_{i,d} \cdot h \cdot s_1)^{\#1} \text{ else} \\ &\quad \text{If } (hd \cdot s_2) \in le + dt \text{ then } \Phi_{i,d} \cdot h \cdot (s_1 \cdot (hd \cdot s_2)) \cdot (tl \cdot s_2) \text{ else} \\ &\quad ((S_{i,d} \cdot h \cdot s_1)^{\#1} \cdot (\Psi(S_{i,d} \cdot h \cdot s_1)^{\#2} \cdot s_2)) \\ \Phi_n \cdot h \cdot s_1 \cdot s_2 &= \text{If } s_2 = () \text{ then } (S_n \cdot h \cdot s_1)^{\#1} \text{ else} \\ &\quad \text{If } (hd \cdot s_2) \in dt \text{ then } \Phi_n \cdot h \cdot (s_1 \cdot (hd \cdot s_2)) \cdot (tl \cdot s_2) \text{ else} \\ &\quad ((S_n \cdot h \cdot s_1)^{\#1} \cdot (\Psi(S_n \cdot h \cdot s_1)^{\#2} \cdot s_2)) \\ \Phi_d \cdot h \cdot s_1 \cdot s_2 &= \text{If } s_2 = () \text{ then } (S_d \cdot h \cdot s_1)^{\#1} \text{ else} \\ &\quad ((S_d \cdot h \cdot s_1)^{\#1} \cdot (\Psi(S_d \cdot h \cdot s_1)^{\#2} \cdot s_2)) \\ \Phi_p \cdot h \cdot s_1 \cdot s_2 &= \text{if } s_2 = () \text{ then } (S_p \cdot h \cdot s_1)^{\#1} \text{ else} \\ &\quad \text{if } (hd \cdot s_2) \in pe \text{ then} \\ &\quad \quad ((S_p \cdot h \cdot s_1 \cdot (hd \cdot s_2))^{\#1} \cdot (\Psi(S_p \cdot h \cdot s_1 \cdot (hd \cdot s_2))^{\#2} \cdot (tl \cdot s_2)) \text{ else} \\ &\quad \quad ((S_p \cdot h \cdot s_1)^{\#1} \cdot (\Psi(S_p \cdot h \cdot s_1)^{\#2} \cdot s_2)) \\ \Phi_e \cdot h \cdot s_1 \cdot s_2 &= \text{if } s_2 = () \text{ then } (S_e \cdot h \cdot s_1)^{\#1} \text{ else} \\ &\quad \text{if } (hd \cdot s_2) = "==" \text{ then} \\ &\quad \quad ((S_e \cdot h \cdot s_1 \cdot (hd \cdot s_2))^{\#1} \cdot (\Psi(S_e \cdot h \cdot s_1 \cdot (hd \cdot s_2))^{\#2} \cdot (tl \cdot s_2)) \text{ else} \\ &\quad \quad ((S_e \cdot h \cdot s_1)^{\#1} \cdot (\Psi(S_e \cdot h \cdot s_1)^{\#2} \cdot s_2)) \end{aligned}$$

The effect of  $\Phi_i$  can intuitively be described as follows: given the beginning of a token in  $L_i$ , then  $\Phi_i$  completes the scanning of this token. This is conceptually much simpler and allows a more efficient implementation than  $\Phi$  in the original definition.

We can now change the definition of  $\Psi$  to utilise the different  $\Phi_i$ :

$$\begin{aligned} \Psi \cdot h \cdot s &= \text{if } s = () \text{ then } () \text{ else} \\ &\quad \text{if } hd \cdot s \in \alpha_i \text{ then } (\Phi_{i,d} \cdot h \cdot (hd \cdot s)) \cdot (tl \cdot s) \text{ else} \\ &\quad \dots \\ &\quad \text{if } hd \cdot s \in \alpha_e \text{ then } (\Phi_e \cdot h \cdot (hd \cdot s)) \cdot (tl \cdot s) \text{ else} \\ &\quad \Psi \cdot h \cdot (tl \cdot s) \end{aligned}$$

Given the simple structure of  $\alpha_i$  the tests in the above definition can be simplified to

$$\begin{aligned} hd \cdot s \in \alpha_d \text{ iff } hd \cdot s &= le \\ \dots \\ hd \cdot s \in \alpha_p \text{ iff } hd \cdot s &= == \end{aligned}$$

#### 2.1.2. Axiomatization of concepts

The revised definitions of  $\Psi$  and  $\Phi_i$  are readily expressed in the rule language

**Language accepted by the verifier.** For example, the definition  
 $\Phi_{id} h s_1 s_2 = \text{if } s_2 = () \text{ then } (S_{id} h s_1)^{*1} \text{ else}$   
 $\text{if } (\text{hd } s_2) \in le + di \text{ then } \Phi_{id} h (s_1, (\text{hd } s_2)) (tl s_2) \text{ else}$   
 $(S_{id} h s_1)^{*1}.(\Psi(S_{id} h s_1)^{*2} s_2)$

is expressed by the rules

**p1i1:** replace  $P\ H\ i\ dent([tab, s1, s2])$  where letter( $\text{hd}(s2)$ ) by  
 $P\ H\ i\ dent([tab, concat(s1, list(\text{hd}(s2))), tl(s2))];$

**p1i2:** replace  $P\ H\ i\ dent([tab, s1, s2])$  where digit( $\text{hd}(s2)$ ) by  
 $P\ H\ i\ dent([tab, concat(s1, list(\text{hd}(s2))), tl(s2))];$

**p1i3:** replace  $P\ H\ i\ dent([tab, s1, s2])$  where  $\neg \text{digit}(\text{hd}(s2)) \wedge \neg \text{letter}(\text{hd}(s2))$  by  
 $\text{concat}([list(S\ ident([tab, s1])), PS([S\ Tr\ ie\ em([tab, s1]), s2]));$

**p1i4:** replace  $P\ H\ i\ dent([tab, s1, null\_sequence])$  by  $list(S\ ident([tab, s1]));$

Similarly we axiomatize other functions of the scanner definition. In addition, we need a theory of sequences. So far we have never made this explicit; rather we have taken properties such as  $(x.y).z = x.(y.z)$  and  $x.y = x.z \Rightarrow y = z$  for granted. A machine checked proof requires formalization of these properties. The complete set of rules required for the scanner verification is included in appendix 3.

## 2.2. Basic algorithm

An implementation of recursively defined functions is straightforward and we could simply implement the functions defining the micro syntax. However, a brief look at the resulting complexity shows that this implementation is not very reasonable: the number of recursive calls is equal to the number of characters in the input.

A more efficient implementation is found by using recursion removal techniques. Consider  $\Psi$  which calls  $\Phi_i$ , which in turn call  $\Psi$  recursively. We will implement  $\Phi_i$  as procedures which append the token scanned to an output file and then return control to  $\Psi$  instead of calling  $\Psi$  recursively.

The scanner operates on the following data structures:

- **outfile** is a file of tokens containing the tokens already scanned.
- **infile** is the input that remains to be scanned.

The procedure **scan**; implementing  $\Phi_i$  is specified as:

$\{ outfile = outfile \wedge infile = infile \wedge hd(infile) = c_0 \wedge c_0 \in L_i \}$

$\{ outfile \Psi(h, infile) = outfile \Phi_i(h, hd(infile), tl(infile)) \}$

where  $h$  is an identifier table.

Given procedures **scan**, as above  $\Psi$  can be implemented as a loop. In this loop **scan**, ( $\Phi_i$ ) will be called repeatedly such that at any point we have:

$$\Psi(h, infile) = outfile \cdot (\Psi(h, infile)).$$

In other words, at any time it is true that scanning the remainder (*input*) with the current identifier table ( $h$ ) and appending this result to the already produced output (*outfile*) is equal to **scan** of the initial input ( $\Psi(h, infile_0)$ ).

The main body of the scanner thus becomes

```
repeat
    if letter(hd(infile)) then scanident else
    if digit(hd(infile)) then scannumber else
    if decimal(hd(infile)) then scandelim else
    if hd(infile) = colon then scancolon else
    if hd(infile) = period then scanper else
    if not eof(infile) then read(infile, c)
    until eof(infile)
invariant
     $\Psi(h_0, infile_0) = outfile \cdot \Psi(h, infile)$ 
```

The correctness of this high level algorithm follows immediately, given the above the definition of **scan**. For example, let us consider a path around the loop where  $hd(infile) \in le$ . It is to prove that

$$\{\Psi(h_0, infile_0) = outfile \cdot \Psi(h, infile) \wedge hd(infile) \in L_{i,id}\}$$

**scanident**  
 $\{\Psi(h_0, infile_0) = outfile \cdot \Psi(h, infile)\}$

The verifier will generate new variable names denoting the final values of variables changed by **scanident**, in this case *infile<sub>1</sub>*, *outfile<sub>1</sub>*, and *h<sub>1</sub>*. The verification condition for this path then becomes

$$\begin{aligned}
 &\Psi(h_0, infile_0) = outfile \cdot \Psi(h, infile) \wedge hd(infile) \in L_{i,id} \\
 &\Rightarrow hd(infile) = c_0 \wedge c_0 \in L_{i,id} \wedge \\
 &\quad (outfile_1 \Psi(h_1, infile_1) = outfile \cdot \Phi_i(h, hd(infile), tl(infile))) \\
 &\Rightarrow \Psi(h_0, infile_0) = outfile \cdot \Psi(h_1, infile_1) \\
 &\Rightarrow \Psi(h_0, infile_0) = outfile \cdot \Psi(h, infile).
 \end{aligned}$$

A little simplification shows that this is equivalent to

$$\begin{aligned}
 &hd(infile) \in L_{i,id} \\
 &\Rightarrow outfile \cdot \Phi_i(h, hd(infile), tl(infile)) = outfile \cdot \Psi(h, infile).
 \end{aligned}$$

This formula is obviously true by the definition of  $\Psi$ . Other cases follow in the same way. Thus, at this point we know that our basic design is sound and, if refined correctly, will lead to a correct program.

Implementation of the procedures `scnr`, poses no problem. Following the definition of  $\Phi$ , the programs are immediate using recursion removal. For the complete program of the scanner see appendix 3.

### 2.3. Implementation details

Now that we decided upon the basic algorithm to be used for the scanner let us look at some of the implementation details.

The abstract program repeatedly refers to the first character of the input stream. To be able to do this in Pascal, we have to read ahead one character and keep the first element of the input stream in a variable. Thus, the remaining input is not given by `in/file` but rather as `(c) in/file` for a variable `c` holding the first element of the character sequence still to be read. However, there is a problem: the input may be empty, `(c) in/file` will never be empty. Tests for the end of file condition have to be recoded. A straightforward solution is to introduce a boolean variable `cvalid` which is true whenever the variable `c` is part of the input string. The actual assertions in the program have been modified accordingly and distinguish the cases `cvalid` and `¬cvalid`. The so revised program is shown in figure 3.

Yet another change can be noticed in figure 3. We implement the identifier table `h` by a data structure `itab` and the representation function `tarep`. Note, that in the formal definition identifier tables are defined as functions  $(\subseteq L, \alpha \rightarrow N) \times (N \rightarrow \{\text{used}, \text{unused}\})$ . In the implementation these functions are represented as data objects. In languages like Lisp where it is possible to dynamically redefine functions (compute on functions) identifier tables could be implemented as functions although this may not be very efficient. We write `apply(h, s)` for function application, i.e. `apply(h, s)` in the program documentation corresponds to `h s` in the formal definition.

`itab` is a record with three components:

- an array of strings which contains identifiers,
- an array of encodings for these identifiers, and
- an integer index pointing to the top element allocated in both arrays.

For the complete definition of this data structure the reader may refer to appendix 3 which contains the complete program text with formal documentation for the scanner.

```

begin
  read(in/file, c);
  cvalid ← true;
repeat
  if letter(c) then scnrident else
  if digit(c) then scnrnumber else
  if delim(c) then scnrnlm else
  if c = colon then scnrcolon else
  if c = period then scnrper else
  if eof(in/file) then cvalid ← false else read(in/file, c)
until not cvalid
invariant
(((cvalid ∧ PS(tarep(itab), in/file)) =
concat(out/file, PS(tarep(itab), concat(is(c), in/file)))) ∧
(¬cvalid ∧ PS(tarep(itab), in/file) = out/file)) ∧ (cvalid ∨ eos(in/file)));
Λisititable(itab);
end

```

Fig. 3

A similar refinement is carried out for the procedures `scnrn`; for the complete program see appendix 3.

### 3. A parser for LS

For the scanner we had a formal definition which was readily transformed into a concrete program. The situation is different in the case of the parser. The parsing function is defined very indirectly, i.e. the result of parsing a program is the unique parse tree if it exists. This parse tree is defined axiomatically; we have to find an algorithm that allows us to compute this tree effectively.

The theory of parsing is well understood and we can resort to existing results in literature [Al74, AE73, AU72, De71, Kn65]. Our compiler uses an LR-parser; the necessary theory of LR-parsing is introduced in the next subsection.

#### 3.1. LR theory

The theory of LR-parsing allows the generation of a *parsing table* which

can be used to drive a parser. Different algorithms for generating tables exist which vary in the size of the generated tables and the class of language that can be handled. For all these methods the parsing algorithms using the tables are identical. Below, we describe the part of LR-theory necessary for the verification of the parser<sup>1</sup>. Although the theory of LR-parsing is well known for some time, a machine checked correctness proof of an LR-parser has not been given.

### 3.1.1. LR-parsing tables

Let  $G$  be a labeled context free grammar satisfying the following conditions:

- There is only one label  $\lambda_0 \in L$  such that  $(P\lambda_0)^{\#1} = s_0$ .
- $P\lambda_0 = \{s_0, \langle n, t \rangle\}$  where  $n \in Nt$  and  $t \in Tm$  such that  $t$  does not occur in any other production.

An LR-parsing table for  $G$  is a mapping  $slr \in St \rightarrow (Nt + Tm) \rightarrow Ac$  where  $Nt$  and  $Tm$  are the set of nonterminals and terminals of  $G$  and  $St$  is a finite set of states with a distinguished initial state  $z \in St$ . The set  $Ac$  is of the form

$$\begin{aligned} Ac = & \{error\} + \\ & \{shift\} \times St + \\ & \{accept, reduce, shift/reduce\} \times L \end{aligned}$$

Let us first give an informal overview showing how these parsing tables are used in the parser. We will then formalize these ideas and show how the correctness of a LR-parser can be proven.

The basic LR-parsing algorithm uses two stacks,  $statestack$ , a stack of states, and  $ntstack$ , a stack of terminals and nonterminals. An input sequence  $input$  contains the string of tokens to be parsed. The algorithm starts with an empty stack  $ntstack$  and with  $statestack$  containing the initial state  $z$ .

The algorithm proceeds by repeatedly performing "parsing actions". Let  $s = top(statestack)$  and  $c = hd(input)$ , then the parsing actions to be performed can be described as follows:

- If  $slr(s, c) = error$  print an error message.
- If  $slr(s, c) = (shift, f)$  push  $f$  on  $statestack$ , push  $c$  on  $ntstack$ , and remove  $c$  from the input.
- If  $slr(s, c) = (reduce, \lambda)$  and  $P\lambda = \langle n, t_1 \dots t_m \rangle$  then remove  $m$  items from  $ntstack$  and  $statestack$  and make  $n$  the first symbol of the input

1.) The tables actually used in the compiler were generated with the SLR generator

- stream (i.e. consider  $\langle n \rangle$  input next).
- If  $slr(s, c) = (shift/reduce, \lambda)$  and  $P\lambda = \langle n, t_1 \dots t_m \rangle$  then remove  $m - 1$  items from  $statestack$  and  $ntstack$ , delete  $c$  from the input, and consider  $n$  the next symbol in the input stream.
- If  $slr(s, c) = (accept, \lambda)$  then the input is parsed correctly.

Successively performing the above steps results in a bottom up parsing of the input. Characters of the input string are shifted on  $ntstack$  until the right hand side of a production  $\lambda$  is on top of  $ntstack$ . Then production  $\lambda$  is applied, that is, the right hand side is replaced by the left hand side. In addition to the above steps the appropriate parse tree has to be constructed; we will discuss this later.

Our informal characterization of the parsing tables is, of course, insufficient for a formal proof. Why should a parser constructed in the above way work? We now present a formal characterization of LR-parsing tables which then enables us to state and prove the correctness of the parser.

With  $W = (Nt + Tm)^*$  we define a relation between sequences of states and sequences of terminals and nonterminals.  $slrrel \subseteq St^* \times W$  as follows:

$$\begin{aligned} slrrel(\langle z \rangle, \langle \rangle) &= \langle \rangle \\ slrrel(u \cdot \langle s_1 \rangle \cdot \langle s_2 \rangle, w \cdot \langle z \rangle) \text{ iff } slr(s_1, w) \wedge slr(s_2, z) &= \langle shft, t, s_2 \rangle \\ slrrel(u, v) \text{ is true, if } u \text{ is of length } n \text{ and if the } i + 1 \text{-st element of } u \text{ is determined by } slr \text{ applied to the } i\text{-th element of } u \text{ and the } i\text{-th element of } v. \text{ In the proof of the parser we show that at any time the relation } slrrel(statestack, ntstack) \text{ holds.} \\ \text{A LR-parsing table has the following relevant properties which enable the construction of a parser. These properties can easily be shown based on the construction algorithm, see for example [De71].} \\ \text{Whenever } slr \text{ calls for a reduce action with production } \lambda \text{ then the right hand side of production } \lambda \text{ is a suffix of } v \\ \text{if } slrrel(u \cdot \langle s \rangle, v) \wedge slr(s, z) = \langle reduce, \lambda \rangle \\ \text{then } v = w \cdot (P\lambda)^{\#2} \text{ for some } w \\ \text{Whenever } slr(s, z) \text{ calls for a shift-reduce action with production } \lambda \text{ then some suffix of } v \text{ followed by } z \text{ constitutes the right hand side of } \lambda \\ \text{if } slrrel(u \cdot \langle s \rangle, v) \wedge slr(s, z) = \langle shift/reduce, \lambda \rangle \\ \text{then } v \cdot z = w \cdot (P\lambda)^{\#2} \text{ for some } w \\ \text{An accept action acts like a shift-reduce action.} \\ \text{If } slrrel(u \cdot \langle s \rangle, v) \wedge slr(s, z) = \langle accept, \lambda \rangle \\ \text{then } (s_0, v \cdot z) \rangle = P \lambda \end{aligned}$$

Whenever  $sir(s, z)$  indicates an error, then  $v$  concatenated with  $z$  is not a prefix of any syntactically valid program.

```
if sirrel(u·s), v) ∧ sir(s, z) = error
  then ∀w.(v·z)·w ∉ L(G))
```

### 3.1.2. The LR-parsing algorithm

We now describe how to use parsing tables introduced above to parse a string and construct a parse tree. The parse tree for an input string is defined in chapter III. It should not be confused with the abstract syntax tree, which is an element of the domain  $\text{asyn}$  of the abstract syntax of  $L_S$ .

Let  $T_T$  be the set of all partial parse trees for  $G$ . We define a relation  $isderiv \subseteq T_T^* \times T_T^*$  which holds for a forest  $r_1, \dots, r_n$  and a string of tokens  $w$  if and only if the leaves of trees  $r_i$  concatenated yield  $w$ . Recursively we define this as:

$$isderiv(\langle \rangle, \langle \rangle)$$

$$isderiv(u, v) \wedge leaves(\emptyset) = w \Rightarrow isderiv(u·\langle \rangle, v·w)$$

Let  $input_0$  be the initial input to the parser, then  $p$  is a parse tree for the string  $input_0$  if

$$isderiv(\langle p \rangle, input_0) \wedge root(p) = s_0$$

where  $s_0$  is the startsymbol of grammar  $G$ . Thus,  $p = \text{parse}(input_0)$  according to the definition of  $\text{parse}$  in chapter II.

In addition to the  $ntstack$  and  $statestack$  the parser uses a third stack  $treetstack$  which is used to construct the parse tree.

The program maintains the following invariant for these stacks:

$$\begin{aligned} & isderiv(treetstack·mktree^*(input), input_0) \wedge \\ & sirrel(statestack, ntstack) \wedge \\ & root(treetstack) = ntstack \end{aligned}$$

Here  $mktree$  constructs a singleton tree:  $mktree = \lambda x.z.T_T$ . The algorithm proceeds by performing appropriate shift and reduce actions until the input string is accepted.

The basic structure of the algorithm is shown in figure 4. The construction of the parse tree is merely indicated by comments.  $treetstack$  is used to build parse trees in the following way: whenever a token is pushed on  $ntstack$ , a singleton tree consisting of just this token is pushed on  $treetstack$ . Then, whenever the parser does a reduction step the partial parse trees corresponding

to the right hand side of the production are on  $treetstack$ . These are used to construct a new tree which is pushed on  $treetstack$ . Similarly we proceed for shift-reduce actions.

Based on the LR-theory above the parsing algorithm can be verified (figure 4 just indicates the principle but is not correct Pascal).

### 3.1.3. Axiomatization

The necessary theories are readily expressed in the verifier's rule language.

For example, the LR-theory is given by

```
%definition of isderiv %
deriv1: infer isderiv(concat(concat(u, list(nd)), v), w) from
  isderiv(concat(concat(u, ts), v), w) ∧ root(ts) = rhs(production(nd));
deriv2: infer isderiv(concat(concat(u, ts), v), w) from
  isderiv(concat(concat(u, ts), v), w) ∧ root(ts) = rhs(production(nd));
deriv3: infer isderiv(concat(concat(rhd(nd), v), w) from
  isderiv(v, w) ∧ null.sequence = rhs(production(nd));
deriv4: infer isderiv(list(mktree(p, n, u, l)), w) from
  isderiv(u, w) ∧ root(u) = rhs(p);

%definition of sirrel %
sirr1: infer sirrel(list(initial_state), null_sequence);
sirr2: infer sirrel(concat(concat(u, list(n)), concat(v, list(n)))) from
  sirrel(u, v) ∧ s = sir(hd(last[1, v]), n).state;
%consequence of 1 and 2 above (requires induction to prove) %
sirr3: infer sirrel(remain(n, s), remain(n, ns));

%properties of LR-parsing tables %
lr1: whenever rhs(sir(s, z).prod)
  from sirrel(st, rt) ∧ sir(s, z).skind = reduce ∧ list(s) = last[1, st];
  infer last[1, length(sir(s, z).prod), ns] = rhs(sir(s, z).prod);
lr2: whenever rhs(sir(s, z).prod)
  from sirrel(st, rt) ∧ sir(s, z).skind = shiftreduce ∧ s = hd(last[1, st])
  infer concat(last[length(sir(s, z).prod) - 1, ns], list(z)) = rhs(sir(s, z).prod);
lr3: from sirrel(st, ns) ∧ sir(hd(last[1, st]), z).syn.skkind = accept
  infer concat(ns, list(z.syn)) = rhs(sir(hd(last[1, st]), z).syn).prod;
len(ns) = 1 ∧
length(sir(hd(last[1, st]), z).syn).prod) = 2 ∧
rhs(sir(hd(last[1, st]), z).syn).prod) = startsymbol ∧
z = eof-symbol;
```

```

begin
  push(statestack, initial_state);
repeat
  act ← alr((top(statestack), hd(input)));
  if act.kind = error then errmsg else
    if act.kind = shift then
      begin
        push(statestack, act.state);
        push(nlistack, hd(input));
        input ← tl(input);
        'build a singleton tree'
      end else
        if act.kind ≠ accept then begin
          if act.kind = reduce then
            begin
              n ← length(act.prod);
              nt ← lh(act.prod);
              nlist ← npop(nlistack, n);
              npop(statetackptr, n);
              'build tree according to the production nt . nlist'
              act ← alr((top(statetackptr), nt));
              end else input ← tl(input);
            while act.kind = shift reduce do
              begin
                n ← length(act.prod);
                'note, only n-1 items are on the stack'
                nlist ← npop(nlistack, n-1);
                npop(statetackptr, n-1);
                nt ← lh(act.prod);
                'build a tree'
                act ← alr((top(statetack), nt));
              end;
              push(nlistack, nt);
              push(statetackptr, act.state);
            end;
            until act.kind = accept
            'accept the input'
          end.

```

### 3.2. Tree transformations

As mentioned earlier the parser will also perform the tree transformations. This section first outlines a slightly modified definition of the tree transformation better suited for an implementation. Next, we show how the transformation step is integrated in the parser.

#### 3.2.1. Building abstract syntax trees

The tree transformation function  $\mathcal{E}$  maps  $WL(Tk)$  into  $Asyn$ . Since  $\mathcal{E}$  is defined recursively, the most natural implementation is to first construct a parse tree and then apply the  $\mathcal{E}$ . However, the particular structure of the tree transformations for  $L_S$  allows a much more efficient implementation. We define a function  $trtr \in WL(Asyn) \rightarrow Asyn$  such that the following holds:

$$\mathcal{E}(\lambda r_1 \dots r_n) = trtr(\lambda \mathcal{E}(r_1) \dots \mathcal{E}(r_n))$$

Of course, it is not always possible to construct such a function  $trtr$  but in our case the structure of  $\mathcal{E}$  is simple enough.

The importance of the definition of  $trtr$  is, that it allows us to compute  $\mathcal{E}(\lambda r_1 \dots r_n)$  without knowing the actual subtrees  $r_i$ . Instead, it is sufficient to know the already transformed subtrees  $\mathcal{E}(r_i)$ . This enables us to construct the abstract syntax bottom up, simultaneously with the construction of the parse tree.

To construct the abstract syntax tree during parsing we add a fourth stack  $aststack$  to the parser. This new stack will be related to  $trestack$  such that at any time we have  $\mathcal{E}^*(trestack) = aststack$ .

We maintain  $aststack$  as follows: whenever a singleton tree  $r$  is pushed onto  $trestack$  we push  $\mathcal{E}(r)$  on  $aststack$ . Consider a reduction step that takes  $r_1, \dots, r_n$  from  $trestack$  and pushes  $(\lambda r_1 \dots r_n)$  on  $trestack$ . For this reduction we remove  $a_1, \dots, a_n$  from  $aststack$  and push  $\mathcal{E}(\lambda a_1 \dots a_n)$  on  $aststack$ . By  $\mathcal{E}^*(trestack) = aststack$  we have  $\mathcal{E}^*(r_1 \dots r_n) = (a_1 \dots a_n)$ . And by definition of  $trtr$  the invariant  $\mathcal{E}^*(trestack) = aststack$  is preserved.

Adding the construction of the abstract syntax to the parser reveals that  $trestack$  as well as  $aststack$  are never really needed except in the formal documentation. Consequently, both these stacks will be virtual data structures that require no implementation.

Fig. 4

### 3.3. Refinement

The actual development of the parser took 10 refinement steps. The intermediate versions of the program can be roughly characterized as follows.

- bare bone version; parsing actions are assumed to be external procedures.
- Stack operations are abstract, so is the input file. There is no output as yet.
- same as previous version but parsing actions are implemented in-line.
- `aststack` is introduced to build an abstract syntax tree.
- a special variable to hold the lookahead is introduced.
- A type `tree` (with pointers) has been defined, `aststack` is implemented as array of trees. However, tree-building operation are still external.
- `statestack` is implemented; now all stacks are implemented, recall that `nstack` and `treetstack` are virtual.
- Reference classes have been introduced together with the necessary documentation to prove that existing trees are unchanged while new ones are being built.
- Auxiliary functions are verified. This step required the change of some entry and exit assertions which in turn called for a reverification of main program with strengthened invariants.
- Routines that build the abstract trees are implemented and verified.

### 3.3.2. Representation

The parser actually does not parse the string contained on the input file but instead parses `input<eof-symbol>` where `eof-symbol` is a new unique token. This token is automatically supplied by the input routine `get` when the input file is empty.

Recall that the grammar has to have the property that  $P\lambda_0 = \{s_0, \{n, t\}\}$  such that  $t$  does not occur in any other production. Thus, given an arbitrary grammar with start symbol  $n$  we will be parsing an augmented grammar with start symbol  $s_0$  and the new production  $\{s_0, \{n, eof-symbol\}\}$  which then is guaranteed to satisfy the condition required in 3.1.1.

The predicate `uniqueof(x)` is true for a sequence  $x$  if `eof-symbol` only occurs as the last element and not somewhere in the middle of  $x$ .

The input sequence of tokens is represented as a pair  $\langle l, \text{input} \rangle$  where  $l$  is the lookahead token and `input` is the remaining input file. The representation function is defined as

$$\text{seqrep}(eof\_symbol, \text{input}) = \langle eof\_symbol \rangle$$

$$\text{seqrep}(l, \text{input}) = \langle l \rangle \cdot \text{input} \cdot \langle eof\_symbol \rangle \text{ where } l \neq eof\_symbol$$

An element of the abstract syntax is represented as a pair of a reference class and pointer. It is based on the principles for representing recursively defined domains (see chapter II), a complete definition of the corresponding representation function `syrep` is contained in appendix 4.

We have to define representations for two stacks. In either case we chose an array together with a starting index and a length. For efficiency reasons we do not use the starting index 1. Since we have to implement the operation `npop` which partitions a stack into two sequences, i.e.

$$\text{npop}(s, n) = \langle s_1, s_2 \rangle$$

such that  $s_1 \cdot s_2 = s$  and  $\text{length}(s_2) = n$  it is convenient to store both sequences in the same array with different starting indices. The representation functions are as follows

$$\begin{aligned} \text{astseqrep}(rc, ar, s, 0) &= \langle \rangle \\ \text{astseqrep}(rc, ar, s, f) &= \text{astseqrep}(rc, ra, s, f - 1) \cdot \langle synrep(rc, ra[s + f - 1]) \rangle \\ \text{stateseqrep}(ar, s, 0) &= \langle \rangle \\ \text{stateseqrep}(ar, s, f) &= \text{stateseqrep}(ar, s, f - 1) \cdot \langle ar[s + f - 1] \rangle \end{aligned}$$

In addition to the standard operations on sequences we introduce the function `last(n, s)` denoting the last  $n$  elements of sequence  $s$ . Similarly, we have the complementary function `remain(n, s)` such that

$$\text{remain}(n, s) \cdot \text{last}(n, s) = s.$$

Yet another point should be mentioned. We are only interested in the case where the program terminates normally. We do not particularly care what happens after an error has been detected except that we require that an error procedure is called and the error be reported to the user. Thus we define a procedure `error`,

```
procedure error;
entry true, exit/else; external;
```

The exit assertion `false` guarantees that the verification condition for any path containing a call to `error` will be true.

### 3.3.3. Reference classes and pointer operations

In chapter II we presented a theorem which simplifies the reasoning about extension operations on pointers. Suitable instances of this theorem can be added to the logical basis of a proof and we will do this in later cases. However, it is also possible to introduce suitable concepts in the assertion language and let the verifier prove necessary instances of the theorem. We demonstrate this in the proof of the parser.

Given a data type  $T$  and a representation function  $\#rep \in T \rightsquigarrow D$ . In chapter II we defined domains represented by a pointer to  $T$  and the associated reference class  $\#T$  as

- $\#rep \in \#T \rightsquigarrow D_P$  for a flat domain  $D_P$  and
- $\#rep \in \#T \rightsquigarrow (D_P \rightarrow D)$ .

The mapping represented by a reference class is partial, i.e. we have

$$\#rep(r) \cdot \#rep(p) = \perp$$

for pointers  $p$  which have not been added to  $r$  by means of the extension operation  $r \cup \{p\}$ .  
We define the predicate *subclass* as follows

$$\text{subclass}(r_1, r_2) \text{ iff } \#rep(r_1) \sqsubseteq \#rep(r_2)$$

i.e. *subclass*( $r_1, r_2$ ) is true if and only if  $r_2$  is an extension of  $r_1$ .

The recursive domain *asyn* of abstract syntax is represented by a pointer and a reference class as outlined in chapter II, we have

$$\text{synrep} \in \#T \times \uparrow T \rightarrow \text{asyn}$$

for a suitable record type  $T$  (for the definition see appendix 4).

Since we require representation functions to be monotonic we have the theorem

$$\text{subclass}(r_1, r_2) \Rightarrow \text{synrep}(r_1, p) \sqsubseteq \text{synrep}(r_2, p)$$

Since the domain *asyn* of abstract syntax is flat, *proper* is defined for *asyn*. Thus, we have

$$\text{subclass}(r_1, r_2) \wedge \text{proper}(\text{synrep}(r_1, p)) \Rightarrow \text{synrep}(r_1, p) = \text{synrep}(r_2, p) \quad (*)$$

In other words, extending a reference class  $r_1$  to  $r_2$  will not change proper objects defined in terms of  $r_1$ .

The predicates *proper* and *subclass* are easily expressed in the verifier's assertion language (see appendix 4). The documentation of the parser includes

suitable assertions involving *proper* and *subclass*. The fact that objects of abstract syntax remain unchanged after extension operations is immediately deducible by the verifier's prover from this documentation and lemma (\*).

### 4. Static semantics

The static semantics is defined with recursive functions  $e$ ,  $t$ , etc (see appendix I). The implementation of these recursive functions is trivial except for one important point. The meaning of recursive type declarations is defined as a least fixed point. This fixed point cannot be computed by simple iteration as this would lead to a nonterminating computation. Therefore, we use the concept of operationalization of fixed points introduced in chapter II.

In this section we discuss the application of operationalization to the definition of types in LS. Further we outline the development cycle of a semantic analysis program in some detail.

#### 4.1. Recursive declaration

In the definition of  $d$  we have the clause

$$d[\![\text{type} \Delta_1; \dots; \Delta_n]\!] \delta = \text{fix}(\text{dt}^*[\![\Delta_1; \dots; \Delta_n]\!]\delta)$$

which describes the effect of a type declaration on the environment. In order to determine the validity of programs we have to be able to compute a representation of the environment given by this fixed point. Note, that it is important here, that we compute the least fixed point. For example, a non-minimal fixed point could be one in which arbitrary new identifiers are defined in the environment. Consequently, weak axiomatization is insufficient.

##### 4.1.1. Operationalization

The solution to the above problem is operationalization of fixed points as introduced in chapter II. The situation here is much more complicated than in the trivial example presented in chapter II. Therefore we will give a slightly simplified description of our solution.

The reader might find the concept of operationalization very unnatural but it has a very intuitive interpretation in our particular case: most compiler constructors would proceed in the very same way. We define the domain  $R$  of unresolved references as:

$$u \in R = (Ty \times Id)^*$$

An element in  $R$  is a list of pairs  $(r, I)$  of a type  $r$  and an identifier  $I$ .

The meaning of such a pair is that the type  $r$  is not yet completely defined. Rather,  $I$  has been used as a type identifier but has not been previously declared;  $r$  is the (as yet unknown) type denoted by  $I$ . This may seem contradictory, how can something unknown be included in a list? The solution is that we store types in a reference class. An “unknown” type is a pointer to an as yet uninitialized element of the reference class. It is this pointer that is included in the list of unresolved references.

Once all type declarations are elaborated, all identifiers in the list of unresolved references are defined (if not, this indicates a semantics error). At this point the unresolved references can be “resolved” by replacing incomplete types by the types denoted by the corresponding identifiers. More precisely, we can update the cells pointed to by pointers in the list of unresolved references.

The key idea of operationalization of a fixed point  $\text{fix } f$  where  $f \in U_s \rightarrow U_s$  is

- to find a representation of the function domain  $U_s \rightarrow U_s$ ,
- to compute the representation of  $f$ , and
- to provide an operation on the representation of  $U_s \rightarrow U_s$  which computes the representation of the least fixed point.

In our case we have  $\text{dt}^* \in T \text{def}^* \rightarrow U_s \rightarrow U_s \rightarrow U_s$  and want to compute  $\text{fix } dt^*[\dots]_S$ , i.e. the argument to  $\text{fix}$  is in  $U_s \rightarrow U_s$ . Suppose we had a representation, say domain  $D$ , for objects in  $U_s \rightarrow U_s$ . Then we want to compute  $\text{dt}^* \in T \text{def}^* \rightarrow U_s \rightarrow D$  and have an operation on  $D$  which computes a representation of the least fixed point.

We will use pairs  $(\varsigma, u) \in U_s \times R$  to represent elements in  $U_s \rightarrow U_s$ . A pair  $(\varsigma, u)$  represents a function, which if applied to  $\varsigma_1$  returns a modified version of  $\varsigma$  in which all unresolved references in  $u$  are resolved by looking up yet undeclared identifiers in  $\varsigma_1$ . To understand the resolution process we have to define the representations used for types, modes, and environments.

A type is represented by a pair consisting of a pointer  $t_p$  and a reference class  $\#t$ . Thus, we have  $\text{tyrep} \in \#T \rightarrow T_p \wedge Ty$ . At this point we are merely interested in functionalities, a precise definition will be given later. Similarly, a mode is represented by a pointer  $m_p$  and a reference class  $\#m$ . But since a component of a mode can be a type a mode also depends of the reference class  $\#t$ ; we have  $\text{modrep} \in \#M \rightarrow \#T \rightarrow M_p \wedge M_d$ .

Finally, environments are represented by records and pointers. Since environments contain modes as well as objects from the abstract syntax the

representation is given by

$$\text{enrep} \in \#A \rightarrow \#M \rightarrow \#T \rightarrow \#E \wedge U_s.$$

If during the evaluation of  $t_f$  we encounter an undefined type identifier, a new cell in  $\#t$  is created but not further initialized. It is only relevant to have a new unique pointer which together with the undeclared identifier is entered in the list of unresolved references. We call those uninitialized types anonymous.

Given this setup resolve is very simple. Suppose we have an element  $f$  in  $U_s \rightarrow U_s$  represented by the pair  $(\varsigma, u)$  where  $\varsigma = \text{enrep}(\#a, \#m, \#t, \#e, z)$ . To determine a representation for  $\text{fix } f$  we simply consider all pairs  $(p, f)$  in  $u$ ; for each pair we determine the type denoted by  $f$  in  $\varsigma$  and update the cell  $\#t \subset p \supset$  accordingly. The analogy to the example presented in chapter II should be obvious.

#### 4.1.2. Revised definition of $t$ and $dt$

We now define revised versions of some valuations which instead of using a second environment to look up forward references build a list of unresolved references whenever they encounter a forward reference during the evaluation. The following arguments are not strictly formal though. To be precise we would have to define the functions below in terms of the representation of types, modes, and environments rather than abstract types, modes, and environments. We trade formality against simplicity, a formal treatment should be obvious though tedious.

$$t_f \in T \text{ typ} \rightarrow U_s \rightarrow R \rightarrow (T_y \times U_s \times R)$$

$t_f[\cdot]_0 u = \text{if } [\text{type } r] : \text{so } [\cdot] \text{ then } (r, s_0, u)$   
 $t_f[I_1, \dots, I_n]_0 u = ([\nu.\text{sub.1}; n], [s_1/I_1], u)$   
 where  $\mu_i = [\text{const}; i; \nu; i]$   
 where  $(\nu, s_2) = \text{newtag } s_0$   
 if  $[\text{const}; i_1, r_1] : \text{so } [E_1]_0$  then  
 $((\text{union } r_1, r_2), s_0, u)$   
 $t_f[\text{array}[T_1]_0 / T_2]_0 u = \text{let } (r_1, s_2, u_2) = t_f[T_1]_0 s_0 u,$   
 $(r_2, s_3, u_3) = t_f[T_2]_0 s_2 u_2 \text{ in}$   
 if  $i \text{ is index } r_1 \text{ then } ([\nu.\text{array}.r_1; r_2], s_4, u_3)$

$t_f[\text{record } I_1; I_2; T_2; \dots; I_n; T_n \text{ end}]_0 u_0 =$   
 if  $i \text{ distinct } (I_1, \dots, I_n)$  then  
 $\text{let } (r_i, s_i, u_i) = t_f[T_i]_0 s_{i-1} u_i \text{ in}$   
 if  $(\nu, s_{i+1}) = \text{newtag } s_i \text{ then }$   
 $([\nu.\text{record}.(I_1, r_1), \dots, I_n, r_n)], s_{n+1}, u_n)$

$t_f[\uparrow I]_0 u = (\{\nu; \uparrow \cdot \tau\}, s_1, u_1)$   
 where  $(\nu, s_1) = \text{newtag } s_0$   
 where  $(\nu, r) = \text{anonymous } [I] u$

Here, **anonymous** creates a new anonymous type which together with the yet undeclared identifier is added to  $u$ . Note, that pointer types are the only case where we allow for forward references to appear; all other types have to be defined previously.

Here is one point where our argument has to be slightly informal because we choose not to present all details of our treatment of recursive declarations. We have to assume that all types  $\tau$  which are returned by **anonymous** are unique and not equal to any existing type. This is apparent if we imagine that in a detailed treatment **anonymous** will not return a type but rather a new pointer to  $\#t$ . Other valuations are changed accordingly:

$$dt_f \in T \text{ def} \rightarrow U_s \rightarrow R \rightarrow (U_s \times R)$$

$$dt_f[I = T]_0 u = \text{let } (r, s_0, u_2) = t_f[T]_0 s_0 u \text{ in } (s_2[[\text{type}; r]/I], u_2)$$

$$dt_f \in T \text{ def}^* \rightarrow U_s \rightarrow R \rightarrow (U_s \times R)$$

$$dt_f[\cdot]_0 u = s_0$$

$$dt_f[\Delta_{i0}; \Delta_{i1}; \dots; \Delta_{in}]_0 u = dt_f[\Delta_{i1}; \dots; \Delta_{in}]_0 s_1 u_1$$
 where  $(s_1, u_1) = dt_f[\Delta_{i0}]_0 u$

#### 4.1.3. Representation of $U_s \rightarrow U_s$

We mentioned that a pair  $(\varsigma, u)$  represents an element in  $U_s \rightarrow U_s$ . Let us now define the necessary representation functions. Not only does a pair  $(\varsigma, u)$  represent an element in  $U_s \rightarrow U_s$ , we also have that a pair  $(r, u)$  represents an element in  $U_s \rightarrow Ty$  and  $(\mu, u)$  represents an element in  $U_s \rightarrow M_d$ . Let us first consider types. We define a function  $\vartheta_{Ty}$  mapping pairs in  $Ty \times R$  into  $U_s \rightarrow Ty$ .

$$\vartheta_{Ty} \in Ty \rightarrow R \rightarrow U_s \rightarrow Ty$$

$$\vartheta_{Ty} \tau u \varsigma =$$

$$\begin{cases} \text{if } u \dots \dots (\tau, J) \dots \dots \text{then } (\text{if } \varsigma[\cdot]_0 : [\text{type}; r_2] \text{ then } r_2) \text{ else} \\ \text{if } r_1 : [\nu.\text{sub.}i; r_2] \text{ then } r \text{ else} \\ \text{if } r_1 : [\nu; \uparrow r_3] \text{ then } [\nu; \uparrow (\vartheta_{Ty} r_3 u \varsigma)] \\ \text{if } r_1 : [\text{nil}] \text{ then } r \\ \text{if } r_1 : [\nu.\text{array}.r_1; r_2] \text{ then } [\nu.\text{array}.(\vartheta_{Ty} r_1 u \varsigma); (\vartheta_{Ty} r_2 u \varsigma)] \\ \text{if } r_1 : [\nu.\text{record}.(I_p, \dots, I_n, p_n)] \\ \text{then } [\nu.\text{record}.(I_p; (\vartheta_{Ty} r_1 u \varsigma), \dots, I_n; (\vartheta_{Ty} r_n u \varsigma))] \end{cases}$$

The intuition behind this definition is simple. If  $r$  is not in the list of unresolved references, the result is  $r$ , regardless of  $\varsigma$ . If however  $r$  is an anonymous type (or if any component type of  $r$  is) then the type returned is the type denoted by  $J$  in  $\varsigma$ , where  $J$  is the identifier associated with  $r$  in the list of unresolved references.

The new valuations  $t_f$  and  $\vartheta_{Ty}$  are related to  $t$  by

$$(\vartheta_{Ty} \tau u \varsigma_1, \varsigma) = t_f[T]_0 s_1$$
 where  $(\tau, \varsigma, u) = t_f[T]_0 s_0$

Similarly we define  $\vartheta_{M_d}$ . If a mode contains an unresolved type, then this type is determined in the given environment:

$\vartheta_{Md} \in Md \rightarrow R \rightarrow (U_s \rightarrow Md)$

$$\begin{aligned}\vartheta_{Md}[\text{var } p] u \varsigma &= [\text{var}:(\vartheta_{Ty} p u \varsigma)] \\ \vartheta_{Md}[\text{vapp } p] u \varsigma &= [\text{vap}:(\vartheta_{Ty} p u \varsigma)] \\ \vartheta_{Md}[\text{valip } p] u \varsigma &= [\text{val}:(\vartheta_{Ty} p u \varsigma)] \\ \vartheta_{Md}[\text{datatype } p] u \varsigma &= [\text{type}:(\vartheta_{Ty} p u \varsigma)] \\ \vartheta_{Md}[\text{constc } p] u \varsigma &= [\text{constc}:(\vartheta_{Ty} p u \varsigma)] \\ \vartheta_{Md}[\text{proc } \mu_1, \dots, \mu_n] u \varsigma &= [\text{proc}:\mu_1, \dots, \mu_n] \\ \vartheta_{Md}[\text{eproc } p] u \varsigma &= [\text{eproc}] \\ \vartheta_{Md}[\text{asfun } \mu_1, \dots, \mu_n] u \varsigma &= [\text{asfun}:\mu_1, \dots, \mu_n, (\vartheta_{Ty} p u \varsigma)] \\ \vartheta_{Md}[\text{dpfun } \mu_1, \dots, \mu_n] u \varsigma &= [\text{dpfun}:\mu_1, \dots, \mu_n, (\vartheta_{Ty} p u \varsigma)]\end{aligned}$$

With the above definitions we can define  $\vartheta_U$  as

$$\begin{aligned}\vartheta_U \varsigma u &= \lambda \varsigma. (\lambda I. \vartheta_{Md}(c^{\#}[I]) u) \varsigma, c^{\#} 2, c^{\#} 3 \\ \text{For } \vartheta_U, dt, \text{ and } dt_f, \text{ we have the relationship:} \\ \vartheta_U \varsigma u \varsigma_1 &= dt[I = T]\varsigma_0 \\ \text{where } (\varsigma, u) &= dt_f[I = T]\varsigma_0 \\ \vartheta_U \varsigma u \varsigma_1 &= dt^*[\Delta]\varsigma_0 \\ \text{where } (\varsigma, u) &= (dt^*[\Delta]\varsigma_0)\end{aligned}$$

#### 4.1.4. Resolving undefined references

Given  $\varsigma$  and  $u$ , a function  $\text{resolve}$  can be defined, such that

$$\text{resolve } \varsigma u = \text{siz}(\vartheta_U \varsigma u)$$

$\text{resolve}$  will compute a representation of the least fixed point by constructing cyclic lists as outlined in chapter II. Clearly,  $\text{resolve}$  cannot be defined in terms of abstract entities like environments, modes, and types alone; we have to

define  $\text{resolve}$  in terms of a particular representation.

```
resolve \varsigma u = if (typerep(\#t, p), J) \u03d5:u then
    let \varsigma' = resolve \varsigma u in
    let envrep(\#a, \#m, \#t, \#e, z) \u03d5: in
        if \varsigma'[J]::[type:typerep(\#t, q)] then
            envrep(\#a, \#m, \#t, \#e, z)
        else
            \varsigma'
```

The proof that  $\text{resolve}$  computes a representation of the least fixed point is analogous to the proof of the general theorem given in chapter II and is not repeated here. Note, that  $\text{resolve}$  can be specified by a weak axiomatization; we will add the definition of  $\text{resolve}$  as well as the definition of  $t_f$  and  $dt_f$  to the logical basis of the proof of the static semantics. An implementation and its correctness is then immediate.

#### 4.2. Development of the program

Given the operationalization of fixed points and the according definitions of  $t_f$  and  $dt_f$  the implementation poses no further problems. In this section we will describe the development process of a part of the static semantics in some detail emphasizing the systematic way in which a program and its specifications can be derived from specifications.

##### 4.2.1. Computing recursive functions

Let us consider the implementation of a procedure that checks expressions for their semantic validity. One clause of the recursive definition of  $e$  is:

```
e[E0][E1]\varsigma = if [v:array r1 of r2]:type e[E0]\varsigma then
    if compatible(r,[type e[E1]]) then
        if isvar e[E0]\varsigma then var:r2] else
            if isval e[E0]\varsigma then [val:r2]
```

As a first step we rewrite this definition in an equivalent first order form which can then be input to the verifier. This translation is straightforward. The only point that requires attention is the type of functions involved. For example  $\text{compatible}$  will not return a boolean result, rather an element in  $\{TT, FF\}^\perp$ . Therefore, truthvalued functions cannot be taken as predicates in the axiomatization. The corresponding rules in the verifier's language are given in figure 5.

Based on this definition we decide to implement  $e$  as a recursive procedure  $C_e$  in our program. The result will be returned in a var parameter  $result$ .

```

rule/file()
constant TT;

replace e(mkindex(E0,E1), zeta) where
  mkarraytype(nu, tau1, tau2) = type(e(E0, zeta)) ∧
  compatible(tau1, type(e(E1, zeta))) = TT ∧
  isvar(e(E0, zeta)) = TT by
  mkvalmode(tau2);

replace e(mkindex(E0, E1), zeta) where
  mkarraytype(nu, tau1, tau2) = type(e(E0, zeta)) ∧
  compatible(tau1, type(e(E1, zeta))) = TT ∧
  isval(e(E0, zeta)) = TT by
  mkvalmode(tau2);

```

Fig. 5

For the first draft we assume that we have functions `isxxx` implementing tests on the abstract syntax as defined earlier. In addition we assume that we have procedures for the decomposition of abstract syntax:

```

function isindex(Ez:Exp):TTFF;
entry true; exit true; external;

```

```

procedure matchindex(Ez:Exp; var E0, E1:Exp);
entry isindex(Ez) = TT; exit Ez = mkindex(E0, E1); external;

```

The true entry and exit specifications cause the verifier to treat `isindex` as a free function symbol. Its meaning is not defined by entry and exit conditions rather it is given as part of the logical basis.

One additional problem is to get the verifier to accept an "abstract" program. In particular this program has to satisfy all declaration and type constraints of Pascal. A simple solution is to define abstract objects as arbitrary enumeration types. We can still take advantage of the system's typechecking but we do not have to define the structure of the data in more detail.

The complete initial program computing `e` (for the case of an indexed expression) is

```

pascal
type Exp = {z1, z2, z3};
U = {z4, z5};

begin
  ...

```

```

Md = (z6, z7);
Typ = (z8, z9);
TTFF = (TT, FF, UU);

procedure error;
entry true; exit false; external;

function isindez(Ez:Exp):TTFF;
entry true; exit true; external;

procedure matchindez(Ez:Exp; var E0, E1:Exp);
entry isindez(Ez) = TT; exit Ez = mkindex(E0, E1); external;

function isarraytype(Ty:Typ):TTFF;
entry true; exit true; external;

procedure matcharraytype(Ty:Typ; var nu:integer; var T0, T1:Typ);
entry isarraytype(Ty) = TT; exit Ty = mkarraytype(nu, T0, T1); external;

function compatible(tau1, tau2:Typ):TTFF;
entry true; exit true; external;

function isvar(mu: Md):TTFF;
entry true; exit true; external;

function isval(mu: Md):TTFF;
entry true; exit true; external;

function mkvalmode(tau:Typ):Md;
entry true; exit true; external;

function type(mu: Md):Typ;
entry true; exit true; external;

procedure C(e:Exp; zeta:U; var result: Md);
entry true; exit e(Ez, zeta) = result;
var E0, E1 :Exp;
tau, tau1, tau2:Typ;
mu, mu1: Md;
nu: integer;
begin
  ...

```

```

if isindex(Ez) = TT then
begin
matchindex(Ez, E0, Ez);
C{E0, zeta, mu};
tau ← typeset(mu);
if not(isarraytype(tau) = TT) then error else
begin
matcharraytype(tau, nu, tau1, tau2);
C{E1, zeta, mu1};
if not(compatible(tau, type(mu)) = TT) then error else
if (isvar(mu) = TT) or (isval(mu) = TT)
then result ← makevalmode(tau2); else error
end else
if ...
else error
end;

```

For this program the verifier produces a set of verification condition. All but one are trivially true since all paths but one contain a call to error. The interesting path is characterized by:

$$\begin{aligned}
& (\text{isindex}(Ez) = \text{tt} \\
\Rightarrow & \text{isindex}(Ez) = \text{tt} \wedge \\
& (Ez = \text{makeindex}(e0\_2, e1\_2) \wedge \\
& e(e0\_2, zeta) = mu\_2 \wedge \\
& \neg(\text{isarraytype}(\text{type}(mu\_2)) = \text{tt}) \\
\Rightarrow & \text{isarraytype}(\text{type}(mu\_2)) = \text{tt} \wedge \\
& (\text{type}(mu\_2) = \text{makearraytype}(nu\_1, tau1\_1, tau2\_1) \wedge \\
& e(e1\_2, zeta) = mu1\_1 \wedge \\
& \neg(\text{compatible}(\text{tau1\_1}, \text{type}(mu1\_1)) = \text{tt}) \wedge \\
& (\text{isvar}(mu\_2) = \text{tt} \vee \\
& \text{isval}(mu\_2) = \text{tt}) \\
\Rightarrow & e(Ez, zeta) = \text{makevalmode}(\text{tau2\_1})) \\
\end{aligned}$$

#### 4.2.2. Refinement

The above verification shows us that the overall design of the program to compute  $\epsilon$  is correct. In order to make  $C\epsilon$  an executable procedure we have to define the representation of the domains  $Exp, Us, Md, Typ$ , and  $T$ . Given a particular representation we can then proceed to refine functions and procedures using these domains.

Let us do this refinement one domain at a time, starting with  $TTFF$ . We let  $truthrep \in \text{boolean}$   $\wedge$   $TTFF$  with

$$\begin{aligned}
truthrep(\text{true}) &= \text{TT} \\
truthrep(\text{false}) &= \text{FF}
\end{aligned}$$

We now refine the program is such a way that the proof we have already given for the abstract version remains valid. The declaration

```

function isindex(Ez; Exp); TTFF;
entry true; exit true; external;
is replaced by the concrete version
function Cisindex(Ez; Exp); boolean;
entry true; exit truthrep(Cisindex) = isindex(Ez); external;
and the call
if isindex(Ez) = TT then ...
becomes
if Cisindex(Ez) then ...

```

We do these changes systematically for all occurrences of  $TTFF$ . The resulting program gives rise to the verification condition given in figure 6.

Applying the definition of  $truthrep$  this verification condition can be transformed into a formula equivalent to the verification condition for the initial program, see figure 7. Consequently, if the refinement is done systematically as outlined above, the verification conditions for the refined version of the program are provable with the same logical basis as the original program. The only additional facts required are the definitions of the representation functions.

#### 4.2.3. Representation

Similar to the implementation of  $TTFF$  we define representations for the abstract domains still undefined and refine the program in the manner outlined above. We will now discuss the representations used for the various domains but will not go through any more refinement steps.

where variables of the form  $z\_1$  are introduced by the verifier. It can easily be seen that this verification condition is provable from our axiomatization of  $\epsilon$  and the verifier's prover handles this case in fractions of a second.

## Static semantics

119

```

(truthrep(cisindex(ex)) = index(ex) ∧
cisindex(ex))
=>
index(ex) = tt ∧
(ex = mkeindex(e0_2, e1_2) ∧
e(e0_2, zeta) = mu_2 ∧
truthrep(cisarraytype(type(mu_2))) = isarraytype(type(mu_2)) ∧
¬cisarraytype(type(mu_2)))
=>
isarraytype(type(mu_2)) = tt ∧
(type(mu_2) ≈ mkearraytype(mu_1, tau2_1) ∧
e(e1_2, zeta) = mu_1 ∧
truthrep(ccompatible(tau1_1, type(mu_1))) =
compatible(tau1_1, type(mu_1)) ∧
¬compatible(tau1_1, type(mu_1)) ∧
truthrep(cisvar(mu_2)) = isvar(mu_2) ∧
truthrep(cisvar(mu_2)) = isval(mu_2) ∧
(cisvar(mu_2) ∨
cisval(mu_2))
=>
e(ex, zeta) = mkevalmode(tau2_1)))

```

Fig. 6

Obviously, the representation of the abstract syntax has to be identical to that used in the parser. A conversion to a different representation is not warranted since the existing one suffices.

Types are represented by the following record structure:

```

type = tnode;
tnode = record
  tkind:(tagrecordtype, tagarraytype,
tagsubtype, tagointerface,
taginttype);
  typeag:integer;
  lwb:integer;
  upb:integer;
  sub1:type;
  sub2:type;
  recotype;
  id:integer
end;
```

The corresponding representation function *typerep* is defined as

## IV. The compiler proof

120

```

(tt = isindex(ex) ∧
isindex(ex))
=>
isindex(ex) = tt ∧
(ex = mkeindex(e0_2, e1_2) ∧
e(e0_2, zeta) = mu_2 ∧
tt = isarraytype(type(mu_2)))
=>
isarraytype(type(mu_2)) = tt ∧
(type(mu_2) = mkearraytype(mu_1, tau1_1, tau2_1) ∧
e(e1_2, zeta) = mu_1 ∧
tt = compatible(tau1_1, type(mu_1)) ∧
tt = isvar(mu_2) ∧
truthrep(cisval(mu_2)) = isval(mu_2)
∨
truthrep(cisvar(mu_2)) = isvar(mu_2) ∧
tt = isval(mu_2))
=>
e(ex, zeta) = mkevalmode(tau2_1)))

```

Fig. 7

```

typerep(#t, p) = if p = nil then ⊥ else
  let r = #t C p D in
    if t.mkind = tagrecordtype then
      [v:record:r] else
      if t.mkind = tagarraytype then
        [v:array:r1:r2] else
        if t.mkind = tagsubtype then
          [v:sub:v1:v2] else
          if t.mkind = tagointerface then
            [v: ↑ r1] else
            [nil]
      where v = typerep(r.typeag),
            where rs = recrep(#t, r.rec),
            where r1 = typerep(#t, r.sust1),
            where r2 = typerep(#t, r.sust2),
            where t1 = intrep(#t, r.lwb),
            where t2 = intrep(#t, r.upb)
    Modes are represented in a similar way:
    mode = ↑ mnode,
    mnode = record
```

modes are represented in a similar way:

*mkind:(tagvalmode, tagupmode, tagfunemode, tagfunsymode,  
 tagvarmode, tagtypemode, tagprocmode, tagprocmode,  
 tagfunconstmode);*

*ty: type;*  
*mlist: mode;*  
*next: mode;*  
*vat: integer*

*end;*

For environments we use a simple minded implementation consisting of two linked lists and an integer number. The two lists represent the functions  $I_d \rightarrow M_d$  and  $N_{\mu} \rightarrow T$  respectively (similar to association lists in Lisp). The number component gives the highest type tag used so far. Since we use type tags consecutively starting with  $\text{int} = 0$ ,  $\text{bool} = 1$  this implicitly gives us a mapping from type tags into  $T$ .

A more sophisticated implementation can be substituted easily without invalidating other parts of the proof.

#### 4.2.4. Auxiliary functions

For abstract syntax, types and modes we provide tests, constructors, and matching functions. Typical examples are:

```
function Cisvalmode(mu: mode): boolean;
global (#tnode, #mnode);
entry true;
exit isvalmode(moderp(#mnode, #tnode, mu)) = truthtable(Cisvalmode);
begin
  Cisvalmode ← mu ↑ .mkind = tagvalmode;
end;

procedure Cisupmode(mul: mode; var mu: mode);
global (#tnode, #mnode, #tnode, #anode);
entry true;
exit moderp(#mnode, #tnode, mu) =
  mkind(moderp(#mnode, #tnode, mu));
begin
  new(mu);
  mu ↑ .mkind ← tagupmode;
  mu ↑ .mlist ← mul;
end;

procedure matchvalmode(mu: mode; var tau: type);

```

```
global (#mnode, #tnode);
entry isvalmode(moderp(#mnode, #tnode, mu)) = "T";
exit moderp(#mnode, #tnode, mu) = mkevalmode(tyerp(#tnode, tau));
begin
  tau ← mu ↑ .ty;
end;
```

A function error is provided and serves the same purpose as in the parser. The false exit condition guarantees that all paths are verifiable if they contain a call to *error*.

No implementation is provided for *o*, *w*, *sf*, and *sp*. An implementation is immediate, given a concrete set of operators and special procedures and functions.

Other auxiliary functions are implemented following their recursive definition. An example is

```
function Cisreturnable(Ty: type): boolean;
global (#tnode, #mnode, #tnode, #anode);
exit truthtable(Cisreturnable) = isreturnable(tyerp(#tnode, Ty));
begin
  Cisreturnable ← Cissubtype(Ty) or
    Cispointertype(Ty) or Cisintype(Ty);
end;
```

#### 4.2.5. The complete program

We have described the development of *Ce*, the implementation of  $\epsilon$ , in some detail. Also, the theory to compute recursive type declarations as been presented. The implementation of the remaining functions is straightforward. In appendix 4 we include several examples, in particular the part of the program dealing with recursive type declarations.

#### 5. Code generation

In all previous parts of the compiler we had to prove that a program computes a particular function. In the code generation the situation is different. The code to be produced does not depend functionally on the input program. We have considerable freedom as to what kind of code we produce subject only to the restriction that it has the same semantics as the input program. One possible approach to code generation is to prove that some abstract "code generating function" produces correct code. This function can then be

taken as specification for a concrete code generator. For example, Milne and Strachey [MS76] provide a suitable code generating function for the language SAL. They give manual proofs of the correctness of this code generation function. Similar proofs are checked mechanically with LCF by A. Cohn [Co79b]. We could have chosen to use a code generating function as specification for our code generation. The techniques for implementation would have been similar to those used in other parts of the compiler. We choose not to do so. Rather, we consider a detailed enough semantics of the source language and an abstract enough definition of the target language such that equivalence between source and target programs can be expressed in the verifier's assertion language.

Still, we need fairly elaborate theorems about source and target semantics and their relation. We omit most of the formal proofs of the necessary semantic theories as these are fairly well understood [MS76]. Rather we concentrate on issues related to program verification; how can the correctness of our code generator be specified? what are the general principles of proof employed?

### 5.1. Principle of code generation

The code generating program is executed after the semantics analysis. Thus, we may assume that the input is a semantically valid abstract syntax tree. The code generation is performed by a set of recursive procedures corresponding to meaning functions of the semantic definition of LS. These functions recursively work their way down the abstract syntax tree and generate the appropriate code. For example let us consider expressions.

The definition of  $\hat{\mathcal{E}}$  contains the following definitional clause for binary operators:

$$\hat{\mathcal{E}}[E_0 \Omega E_1] \rho \gamma = \mathcal{R}[E_0] \rho; R[E_1] \rho; binop[\Omega] \gamma$$

Corresponding to  $\hat{\mathcal{E}}$  and  $\mathcal{R}$  we have a procedure *AEcode* and *Recode*. In the case of binary operators *AEcode* proceeds as follows:

- Call *Rcode* for the first expression. The resulting code will guarantee that after its execution the value of the first expression is on top of the stack.
- Call *Rcode* for the second expression; append the new code to the previously generated code. Now, at this point during runtime the results of both expressions will be the two topmost elements of the stack.
- Generate code to execute the binary operation in question. This code will remove the two topmost elements and push the result of the operation.

Thus, altogether the code  $\hat{\mathcal{E}}$  generates is what leaves the result of the expression  $E_0 \Omega E_1$  on top of the stack.

*Rcode*: “self is similar to *Acode*; it merely adds possible coercions of L-values ( $L(V)$ ) to  $V$ , i.e. ( $V$ )

One might ask „What does ‘self’ mean?“ Well, we have to generate code that is “self” in the sense that it is “self” in the sense that it is associated with this variable. But since  $E_0$  and  $E_1$  are in the target language are different we have to find the “self” in the target language. And the necessary information about locations assigned to variables in the target language we introduce a “compile time environment”. A “compile time environment” is similar to an environment except that it maps identifiers into locations in the target semantics rather than into locations in the source semantics.

If the code generating procedures encounter a declaration new locations on the target machine are allocated. If the corresponding identifiers are bound to these new locations in the compile time environment

By the choice of the definitional methods for source and target languages addressing is relative in *LS* and *LT*. In the produced code the relative addresses are hard-wired while the base addresses are only known at run time, the final address computation determining the absolute address is done at execution time.

Special considerations are necessary to treat labels properly. Consider a conditional statement *If E then  $\Gamma_0$  else  $\Gamma_1$*  for which we generate the following code:

- generate two new unique labels  $N_0$ , and  $N_1$ .
- generate code to evaluate *E*
- generate a conditional jump to  $N_0$  if *E* is false
- generate code for  $\Gamma_0$
- generate an unconditional jump to  $N_1$
- define label  $N_0$  in the output code
- generate code for  $\Gamma_1$
- define label  $N_1$  in the output code

This code sequence introduces new labels in the target program which have no counter part in the source language. Recall that the definition of *LT* considers all labels of a target program at the same time (outer level); i.e. one fixed point determines all label continuations. The equivalence of a target program with a source program is therefore very hard to establish if there is

not a one to one correspondence of labels.

We will solve this problem in two steps. First, we change the definition of the semantics of conditionals in the source language such that it uses explicit labels and jumps. We prove that this new definition is equivalent to the original one.

In addition we introduce a block construct in the target language that allows us to encapsulate labels. A block  $(I_1, \dots, I_n)$  is considered a single instruction, internally a block consists of a list of instructions, possibly other blocks. All labels defined inside a block are local to this block. Thus, for the above example we will generate one block of code. Auxiliary labels are hidden inside the block, furthermore we have a one to one correspondence of statements of  $LS$  and instructions (blocks) of  $LT$ .

### 5.2. Modified semantics definitions

In this section we define a block construct as an extension to the target language  $LT$ . We show how the extended language relates to the original language. We prove a theorem stating that under certain weak conditions programs in the extended language can be transformed into programs in  $LT$ .

Furthermore, we present a revised definition for some statements of  $LS$  in terms of conditionals and explicit jumps. We prove the equivalence with the original definition.

#### 5.2.1. A structured target language

To solve the label problem mentioned above we will introduce a block structure in the target language in such a way that label values within a block can be determined independently from the rest of the program. This will allow to translate, say a while loop, into some code block whose correctness we can determine without looking at the context.

We will prove a theorem stating that under very weak assumptions a block of code can be textually expanded (block boundaries can be omitted) without altering the meaning of the program. In generating code in the compiler we will prove that the necessary conditions are satisfied. Therefore the generated block structured code will be equivalent to the same code with block boundaries ignored.

Let us define a different language  $LT2$  by adding the instruction  $(I_1, \dots, I_n)$  to the instruction set of  $LT$  with the semantics

$$\mathcal{M}[(I_1, \dots, I_n)]\rho\gamma = \text{if } \text{distinct}(I_1, \dots, I_n) \text{ then } \mathcal{M}^*(I_1, \dots, I_n)\rho\gamma$$

where  $\rho_1 = \text{fix}(\lambda p. \rho[L[I_1, \dots, I_n]\rho]/[I_1, \dots, I_n])$ .

In effect,  $LT2$  allows for programs as instructions. The importance is that labels within each program are only visible locally. We will now show that a program in  $p \in LT2$  is equivalent to  $p$  with all parentheses omitted, provided all labels in the resulting program are distinct.

Let  $S_i$  be sequences of statements and  $p, q$ , and  $r$  programs such that

$$p = S_1 q S_3, \quad q = (S_2), \quad \text{and } r = S_1 S_2 S_3$$

We call  $r$  the "flat" version of  $p$ , i.e. in the flat version block boundaries are omitted. Let us introduce the abbreviations  $L_i = L[S_i]$ ,  $I_i = I[S_i]$ , and  $M_i = M[S_i]$ .

The condition under which a program is equivalent to its flat version is that labels defined in  $S_1$  and  $S_3$  are distinct from those defined in  $S_2$  and that there are no jumps into inner blocks. Clearly, if a  $LT$  program is well formed, jumps into inner blocks are impossible since inner labels are not visible.

We write  $\rho =_2 \bar{\rho}$  for two environments that only differ in the values bound to labels in  $I_2$ . Formally,  $\rho =_2 \bar{\rho}$  holds if  $\forall l. \notin I_2 \rightarrow \rho[l] = \bar{\rho}[l]$ . The  $LT$  program  $p$  (as above) is well formed if  $\rho =_2 \bar{\rho}$  implies that  $M_1\rho = M_1\bar{\rho}$  and  $M_3\rho = M_3\bar{\rho}$ .

**Theorem** If  $I_1, I_2, I_3$  are distinct and if  $p$  is well formed then  $\mathcal{M}[p]\rho\gamma = \mathcal{M}[r]\rho\gamma$ .

**Proof** Let  $\rho_p, \rho_q$ , and  $\rho_r$  be the environments in which the (top level) statements of programs  $p, q$ , and  $r$  are evaluated. By definition of  $\mathcal{M}$  we have

$$\begin{aligned} \rho_q &= \text{fix } \lambda p. \rho [L[\bar{\rho}; M_3\bar{\rho}\gamma / I_2]] \\ \rho_p &= \text{let } \rho_{\bar{p}} = \text{fix } \lambda \bar{\rho}. \bar{\rho} [L_1[\bar{\rho}; M_3\bar{\rho}\gamma / I_2]] \text{ in} \\ &\quad \text{fix } \lambda \bar{\rho}. \rho [L_1[\bar{\rho}; M_2\bar{\rho}\gamma / I_1]] [L_2[\bar{\rho}\gamma / I_3]] \\ \rho_r &= \text{fix } \lambda \bar{\rho}. \rho [L_1[\bar{\rho}; M_2\bar{\rho}; M_3\bar{\rho}\gamma / I_1]] [L_2[\bar{\rho}; M_3\bar{\rho}\gamma / I_2]] [L_3[\bar{\rho}\gamma / I_3]] \end{aligned}$$

First we prove that  $\rho_q = \rho_r$ ; that is, the instruction sequence  $S_2$  is evaluated in the same environment as the flat version  $r$  of  $p$ . Let us rewrite the fixed points as equivalent recursion equations. We can define  $\rho_p$  and  $\rho_q$  by a mutually recursive definition as follows

$$\begin{pmatrix} \rho_q \\ \rho_p \end{pmatrix} = \begin{pmatrix} \rho_p [L_2[\bar{\rho}; M_3\bar{\rho}\gamma / I_2]] [L_3[\bar{\rho}\gamma / I_3]] \\ \rho_p [L_1[\bar{\rho}; M_2\bar{\rho}; M_3\bar{\rho}\gamma / I_1]] [L_2[\bar{\rho}; M_3\bar{\rho}\gamma / I_2]] \end{pmatrix}$$

Similarly we can define  $\rho_p$ , as

$$\begin{pmatrix} \rho_r \\ \rho_y \end{pmatrix} = \begin{pmatrix} \rho_p [L_2 \rho_p; M_3 \rho_p \gamma / l_2] \\ \rho [L_1 \rho_p; M_3 \rho_p \gamma / l_1] [L_3 \rho_p \gamma / l_3] \end{pmatrix}$$

First, observe, that the recursion equation for  $\rho_q$  and  $\rho_z$  are identical (up to renaming of the variables). Thus, we can eliminate  $\rho_z$  without changing the fixed points for  $\rho_q$  and  $\rho_p$ .

$$\begin{pmatrix} \rho_q \\ \rho_p \end{pmatrix} = \begin{pmatrix} \rho_p [L_2 \rho_q; M_3 \rho_q \gamma / l_2] \\ \rho [L_1 \rho_q; M_2 \rho_q; M_3 \rho_q \gamma / l_1] [L_3 \rho_q \gamma / l_3] \end{pmatrix}$$

Since we have

$$\rho_r = \rho_y [L_2 \rho_r; M_3 \rho_r \gamma / l_2]$$

$\rho_r$  and  $\rho_y$  only differ for labels in  $l_2$ . Thus, by definition of  $\equiv_2$  we have  $\rho_r \equiv_2 \rho_y$ . By our assumptions of well formedness we have  $L_1 \rho_r = L_1 \rho_y$ ,  $L_3 \rho_s = L_3 \rho_y$ , and  $M_3 \rho_s = M_3 \rho_y$ . Therefore, we can rewrite the definition of  $\rho_s$  as

$$\begin{pmatrix} \rho_r \\ \rho_y \end{pmatrix} = \begin{pmatrix} \rho_y [L_2 \rho_p; M_3 \rho_p \gamma / l_2] \\ \rho [L_1 \rho_y; M_3 \rho_y \gamma / l_1] [L_3 \rho_y \gamma / l_3] \end{pmatrix}$$

But now  $\rho_r$ ,  $\rho_y$  and  $\rho_p$  are defined by isomorphic recursion equations and thus must have equal fixed points. We conclude  $\rho_p \equiv \rho_y$  and  $\rho_q \equiv \rho_r$ , and  $\rho_p \equiv_2 \rho_r$ .

The theorem follows immediately. Since  $\rho_p \equiv_2 \rho_r$ , we have  $M[S] \rho_p = M[S] \rho_r$ , and  $M[S_3] \rho_p = M[S_3] \rho_r$ . Furthermore,  $M[g] \rho_p = M[S_2] \rho_p = M[S_2] \rho_r$ , which completes the proof. ■

**Corollary** If the restrictions on the use of labels for the theorem are satisfied, then every program is equivalent to its flat version

**Proof** Immediate, since any program has only a finite nesting the theorem can be applied to the program repeatedly. ■

### 5.2.2. A modified definition of LS

We add the statement `unless E goto N` to the abstract syntax of LS. We define

$$B[\text{unless } E \text{ goto } N] \rho = R[E] \rho; \text{Cond}(\gamma, \rho[N])$$

This new statement is not available to the LS programmer, i.e. no external syntax is provided. Rather, the `unless` statement is used to define conditionals and loops as follows:

$$\begin{aligned} B[\text{if } E \text{ then } \Gamma_0 \text{ else } \Gamma_1] \rho &= C[\text{unless } E \text{ goto } N_0; \\ &\quad \text{begin } \Gamma \text{ end;} \\ &\quad \text{goto } N_1; \\ &\quad \text{No:begin } \Gamma_1 \text{ end;} \\ &\quad N_1:\text{dummy}] \rho \gamma \\ &\quad \text{where } \varsigma[N_1] = \text{false}, \varsigma[N_0] = \text{false} \wedge N_0 \neq N_1 \\ B[\text{while } E \text{ do } \Gamma \text{ od}] \rho \gamma &= C[\text{if } N_0: \text{unless } E \text{ goto } N_1; \\ &\quad \text{begin } \Gamma \text{ end;} \\ &\quad \text{goto } N_0; \\ &\quad N_1:\text{dummy}] \rho \gamma \\ &\quad \text{where } \varsigma[N_1] = \text{false}, \varsigma[N_0] = \text{false} \wedge N_0 \neq N_1 \\ B[\text{repeat } \Gamma \text{ until } E] \rho \gamma &= C[\text{if } N_1: \text{begin } \Gamma \text{ end;} \\ &\quad \text{unless } E \text{ goto } N_1] \rho \gamma \\ &\quad \text{where } \varsigma[N_1] = \text{false} \end{aligned}$$

**Theorem** The new definitions for conditionals, while and repeat loops are equivalent to those of the original definition (appendix 1).

**Proof** We consider only while loops. Other cases follow by similar arguments. By the original definition we have

$$B[\text{while } E \text{ do } \Gamma \text{ od}] \rho \gamma = \text{fix } (\lambda \dot{\gamma}. \mathcal{R}[E] \rho; \text{cond}(C[\Gamma], \rho, \gamma))$$

It is to show that

$$\begin{aligned} \text{fix } (\lambda \dot{\gamma}. \mathcal{R}[E] \rho; \text{cond}(C[\Gamma], \rho, \gamma)) &= \\ C[\text{if } N_0: \text{unless } E \text{ goto } N_1; \\ &\quad \text{begin } \Gamma \text{ end;} \\ &\quad \text{goto } N_0; \\ &\quad N_1:\text{dummy}] \rho \gamma \\ &\quad \text{where } \varsigma[N_1] = \text{false}, \varsigma[N_0] = \text{false}, N_0 \neq N_1, \end{aligned}$$

Let us abbreviate  $N_0: \text{unless } E \text{ goto } N_1$  and  $\Gamma$  the right hand side of the above equation as  $\Gamma_{new}$ . By the definition of  $C$  the right hand side of the above equation becomes

$$C[\Gamma_{new}] \rho$$

where  $\delta = \{[\text{true}, \text{true}]/(N_0, N_1)\}$  and  $\rho = \text{fix } (\lambda p[\lambda^{\{J\}} [\Gamma_{n,w}] \{p/\{N_0, N_1\}\})$ . Evaluating  $J$  we get  $\rho = \text{fix } (\lambda p[\lambda^{\{C[\Gamma_{n,w}] \{p/\{N_0, N_1\}\}}])$ .

Now consider  $C[\Gamma_{n,w}] \{p\}$ ; after some simplification we get

$$C[\Gamma_{n,w}] \{p\} = R[E] \{p\} \text{cond}(C[\Gamma] \{p\}, \bar{p}[N_0]).$$

By the general rule

$$\text{fix}(\lambda p. \rho[\{p/z\}/z]) = \rho[\text{fix } (\lambda \gamma. J(\gamma)) / z]$$

we derive

$$\rho[N_0] = \text{fix } (\lambda \gamma. R[E] \{p\} \text{cond}(C[\Gamma] \{p\}, \bar{p}\gamma))$$

Since  $\rho[N_i] = \text{false}$  the new labels  $N_0$  and  $N_1$  cannot appear as global labels of either  $E$  or  $\Gamma$ . Therefore we have

$$C[\Gamma] \{p\} = C[\Gamma] \{p\}$$

and

$$R[E] \{p\} = R[E] \{p\}$$

which proves the theorem.  $\blacksquare$

### 5.3. Relation between $LS$ and $LT$

Asserting that a  $LS$  and a  $LT$  program have the same meaning is easy since their respective domains of meanings are identical ( $N^* \rightarrow N^* + ER$ ). This situation is more complicated at intermediate stages of a computation.

Consider again the example of expressions given earlier. Not only do we have to prove that the correct result is pushed on the stack by the generated code; we also have to prove that evaluating an expression causes "corresponding side effects". In the following discussion we will make the notion of "corresponding meaning" precise.

We introduce a set of predicates that relate objects in the source language  $LS$  to corresponding objects in the target language. For example, answers ( $A$ ) in  $LS$  are related to answers in  $LT$  by equality, i.e. answers correspond if and only if they are equal. We will use indices  $S$  and  $T$  to disambiguate objects of source and target semantics that are denoted by the same symbol. We define

$$\delta_A(a_S, a_T) \equiv a_S = a_T$$

For other domains  $D$  the relations  $\delta_D$  are not as trivial.

The central question is how objects of  $LS$  are stored in the target language and how the relation between locations of  $LS$  and  $LT$  can be expressed. Also, it has to be clarified how other declared objects such as procedures, functions and labels are implemented in  $LT$ .

We use "compile time environments" to relate labels and procedure identifiers of  $LS$  to labels of  $LT$  and variable identifiers to relative addresses of  $LT$ .

"Relative storage maps" relate relative locations of  $LS$  to relative locations of  $LT$ . Finally, "absolute storage maps" relate absolute locations of  $LS$  and  $LT$ .

Given a variable, procedure, or function identifier or a numeral denoting a label in  $LS$ . A compile time environment  $y \in Y$  specifies what the corresponding objects in  $LT$  are. We define  $Y$  as

$$y \in Y = (Id \rightarrow N) \times (Id \rightarrow N) \times (Num \rightarrow N) \times N$$

The first component is a mapping from identifiers to relative locations; the second component maps procedure and function identifiers into labels. The component  $Num \rightarrow N$  maps labels of  $LS$  into labels of  $LT$ . The last component  $N$  is mainly used for bookkeeping purposes; it specifies the number of memory cells that are allocated in the current lexical level.

### 5.3.2. Storage allocation

Let us now decide how data objects of  $LS$  are to be stored in  $LT$ . All data objects other than arrays and records are stored in one cell. All elements of an array have the same size (require the same number of cells); furthermore the element size is known at compile time. Therefore, all elements of an array are stored in consecutive locations. The first of these locations is the address of the array. The situation is similar for records. Since number and size of the components are known at compile time all elements are stored consecutively. Again, the first address is the address of the record.

To access a component of a record or an array we need the address (first element) of the record or array. Given an index or a selector the offset from the first element address can be computed.

Let  $\text{size}$  be a function determining the number of cells required to store a variable of type  $r$ .

```

size ∈ Ty → N

size[v; sub; i1;i2] = 1
size[v; t] = 1
size[init] = 1
size[v; array; r1;r2] = If r1:[v; sub];i1;i2] then
  (size r2)*(i2 - i1 + 1)
size[v; record; l1;r1, ..., ln;rn] = (size r1) + ... + (size rn)

```

Given an object of type  $r$  with the  $L$ -value  $a$  which is assigned the address  $z$  in the target language we can define how addresses of components of this object are related to target addresses. We define a function  $ad \in Lv_s \rightarrow Ty \rightarrow N \rightarrow ((L_s \times N) \rightarrow T)$  which given a  $L$ -value, a type and an address in  $LT$  returns a relation  $\in (L_s \times N) \rightarrow T$  between locations of  $LS$  and addresses in  $LT$ .

```

ad a r z = If a ∈ L, then {a, z} else
  if r:[v; array; [v; sub; i1;i2]; r2] then
    ad(a[i1])r2 z ∪
    ad(a[i1 + 1])r2(z + size r2) ∪
    ...
    ad(a[r2])r2(z + (i2 - i1)*size r2) else
    If r:[v; record; l1;r1, ..., ln;rn] then
      ad(a[l1])r1 z ∪
      ad(a[l2])r2(z + size r1) ∪
      ...

```

### 5.3.3. Storage maps

We define the domain of storage maps as

$$s \in \Sigma = (L_r \times N) \rightarrow T.$$

A pair  $(\alpha, m)$  is in a storage map  $s$  if the relative location  $\alpha$  in  $LS$  corresponds to the location  $m$  in  $LT$ . For each scope a different storage map obtains.

Note, that a storage map is only defined for relative locations  $(L_r)$ , e.g. nothing is said about array objects in  $I \rightarrow Lv_s$ . This is justified since complex  $L$ -values  $\in I \rightarrow Lv$  are not stored in the target program; only elements

2.) To be precise, relations should be defined as mappings into  $(TT, \perp)$  with  $\perp \subseteq TT$ . This is important to guarantee the existence of recursively defined predicates. For details see [MS76,Re74].

of arrays are stored. If  $LS$  had dynamic arrays, array descriptors would be necessary for an implementation. In this case the descriptor could be viewed as representing a complex  $L$ -value.

For given environments  $\varsigma$  and  $\rho$  and a compile time environment  $\gamma$  the relative storage map for the current scope is given by  $\Phi_\Sigma \in U_\varsigma \rightarrow U \rightarrow Y \rightarrow \Sigma$  as

$$\Phi_\Sigma \rho \gamma = \{(x, z) \mid \exists J \in Id, \rho^{\#^3}[J] = level \rho \wedge (x, z) \in ad(\rho^{\#^1}[J])(type, s[I](\nu[J]))\}$$

Absolute storage maps are defined as  $l \in \Lambda = L \times N \rightarrow T$ ; they relate absolute addresses of  $LS$  and  $LT$ .

For given  $s \in \Sigma$  and a source and a target display the corresponding absolute storage map  $l \in \Lambda$  is determined by the function  $\Phi_\Lambda \in \Sigma \rightarrow X \rightarrow D \rightarrow \Lambda$  as

$$\Phi_\Lambda s \chi \delta = \{(x^{\#2} n z, \delta^{\#2} n + y) \mid n = \delta^{\#4} \wedge (x, y) \in s\}$$

Relations given by  $\Phi_\Lambda$  only relate static locations ( $\in L_s$ ). We assume that  $L_d$  and  $N$  are isomorphic, thus a relation  $l_h \in \Lambda$  can be established once and for all.

The absolute storage map will vary during program execution. For example, exiting a procedure will change the locations that are relevant for the rest of program execution. We define a structure that is similar to display<sup>1</sup> which allows us to keep track of varying storage relations at all times.

A run time storage relation ( $\in Q$ ) is defined as

$$\beta \in Q = (N \rightarrow Q) \times Q \times \Lambda.$$

For the dynamic and all static predecessors of the current display  $\beta$  defines a run time storage relation. The component  $\Lambda$  is the dynamic storage map that obtained at the call of the current procedure.

Entering a new procedure  $\beta$  has to be incremented. Such a change of  $\beta$  is described by  $up_Q \in Q \rightarrow X \rightarrow D \rightarrow N \rightarrow \Sigma$  as

$$up_Q \chi \delta \beta n s = \begin{cases} \text{fix } \lambda \beta.(f, \beta, l) \\ \text{where } f = \lambda \nu. \text{ If } \nu < n \text{ then } \beta^{\#1}, \\ \quad \text{else if } \nu = n \text{ then } \beta \text{ else } \perp \\ \text{where } l = \Phi_\Lambda \chi \delta s \bigcup \beta^{\#3} \end{cases}$$

Here,  $s \in \Sigma$  is the storage map that obtains at the point of the procedure call;  $\chi$  and  $\delta$  are corresponding source and target displays and  $n$  is the lexical level of the procedure.

### 5.3.4. Relations between domains

#### 5.3.4.1. Values

We use  $\vartheta_V$  to relate values of  $LS$  and  $LT$ . Values in  $LT$  are in  $N$ ; values of  $LS$  are either locations or index values ( $I \subseteq N$ ). Since locations are related by absolute storage maps,  $l \in \Lambda$  has to be argument to  $\vartheta_V$ ; we define

$$\vartheta_V(\epsilon, l, n) \equiv \text{if } \epsilon \in I \text{ then } \epsilon = n \text{ else } l(\epsilon, n)$$

#### 5.3.4.2. Stores

Depending on a particular  $l$  we can define when a store in  $LS$  is equivalent to a memory state in  $LT$ . A store  $\sigma$  is equivalent with a memory  $\mu$  if input and output files are equal, i.e.  $\sigma^{#2} \wedge \sigma^{#3} = \mu^{#3}$ , and if corresponding locations are mapped into corresponding values, i.e.  $\forall a, n. l(a, n) \Rightarrow \vartheta_V(\sigma a, l, \mu n)$ . This is not quite correct though, since we are only interested in used locations. The last component of  $\mu$  specifies how much memory is used in the heap; therefore we can write  $\forall a, n. l(a, n) \wedge n > \mu^{#4} \Rightarrow \vartheta_V(\sigma a, l, \mu n)$ . In addition all  $a \in L_d$  corresponding to  $n \leq \mu^{#4}$  have to be unused. We require:  $\forall a, n. l(a, n) \wedge n \leq \mu^{#4} \Rightarrow \sigma a = \text{unused}$ . Altogether we define

$$\begin{aligned} \vartheta_S(\sigma, l, \mu) \equiv & \sigma^{#2} = \mu^{#2} \wedge \\ & \sigma^{#3} = \mu^{#3} \wedge \\ & \forall a, n. l(a, n) \wedge n > \mu^{#4} \Rightarrow \vartheta_V(\sigma a, l, \mu n) \wedge \\ & \forall a, n. l(a, n) \wedge n \leq \mu^{#4} \Rightarrow \sigma a = \text{unused} \end{aligned}$$

#### 5.3.4.3. Displays

Two displays  $X$  and  $\delta$  and a run time storage relation  $\beta$  are related if the dynamic and all static predecessors are related, if the current lexical levels are equal, and if the base address in  $\delta$  is greater than all memory used in the absolute storage map in  $\beta$ . Displays of  $LT$  have two additional components which have no counterparts in displays of  $LS$ . The relation of these components to the source language will become apparent if we consider procedure values.

$$\begin{aligned} \vartheta_X(X, \beta, \delta) \equiv & \theta_X(X^{#3}, \beta^{#2}, \delta^{#3}) \wedge X^{#4} = \delta^{#4} \wedge \\ & \forall n < \delta^{#4}. \vartheta_X(X^{#1} n, \delta^{#1} n) \wedge \\ & \text{used}(\beta^{#3}) < \delta^{#2} \delta^{#4} \wedge \end{aligned}$$

Here  $\text{used}$  determines the greatest memory location in  $LT$  that is in use according to  $l \in \Lambda$ . We define

$$\text{used}(l) = \max\{n \mid (a, n) \in l\}$$

If a procedure is entered and  $X, \beta$  and  $\delta$  are incremented accordingly, the relation  $\vartheta_X$  is preserved. We have the lemma

$$\vartheta_X(X, \beta, \delta) \Rightarrow \vartheta_X(\text{ups} X n, \text{up} Q X \delta \beta n, \text{up} r \delta \gamma n m k)$$

The proof is immediate by the definition of  $\text{ups}$ ,  $\text{up} Q$ , and  $\text{up} r$ .

#### 5.3.4.4. Continuations

The domain in  $LT$  corresponding to generalized dynamic continuations  $G_d$  is  $S \rightarrow M \rightarrow A$ . We define

$$\begin{aligned} \vartheta_{G_d}(\gamma, l, \theta) \equiv & \text{if } \gamma \in C \text{ then } \forall \sigma, \mu. \vartheta_S(\sigma, l, \mu) \Rightarrow \vartheta_A(\gamma \sigma, \theta(l, \mu)) \\ & \text{else } \forall \epsilon, n. \vartheta_V(\epsilon, l, n) \Rightarrow \vartheta_G(\gamma \epsilon, l, \lambda \sigma. \theta(l, \sigma)) \end{aligned}$$

Given a compile time environment  $\gamma \in Y$  we can define when two continuations correspond:

$$\vartheta_G(\gamma_S, \epsilon, \gamma_T) \equiv \forall \chi, \beta, \delta. \vartheta_X(\chi, \beta, \delta) \Rightarrow \vartheta_G(\gamma_S X, l, \gamma_T \delta)$$

where  $l = \Phi_{\Lambda} \sigma X \delta \cup \beta^{#3} \cup l_0$ . That is,  $l$  is the relation of all dynamic locations ( $l_h$ ) and all used static locations.

#### 5.3.4.5. Procedures and functions

There is no analog to procedure or function values in  $LT$ . Procedure and function values correspond to continuations in  $LT$  in the following way.  
If a label  $j$  in a  $LT$  program marks the beginning of the code for a procedure, then this procedure can be called by the instruction

$$CALL j \ n \ m \ k$$

where  $n$  is the lexical level of the procedure,  $m$  is the number of parameters, and  $k$  is the number of memory cells allocated in the current environment.

Correspondence between a procedure identifier  $p$  and the label  $j = \gamma[p]$  is given if a procedure call to  $j$  has the same effect as executing the procedure  $p$ . Formally we define

$$\begin{aligned} \vartheta_p(n, m, \gamma) \equiv & \forall \epsilon, \gamma_S, \gamma_T. \vartheta_G(\gamma_S, \gamma_T) \Rightarrow \vartheta_G(\pi \gamma_S, \epsilon, \gamma_T) \\ \text{where } \gamma &= \lambda \delta \sigma. \gamma(\text{up} \delta \gamma_T(n + 1)(\text{length } \sigma - m), k) \sigma \\ \text{where } k &= \text{used}(\epsilon) \end{aligned}$$

It is important, that the value  $k$  is determined from the storage map at call time since it is used to allocate (yet unused) memory for the activation record of the current call.

#### 5.3.4.6. Environments

The only objects bound in environments of  $L_T$  are labels; there are no identifiers in  $L_T$ . For  $\rho_S$  in  $L_S$  to be equivalent to  $\rho_T$  in  $L_T$  the following conditions must hold.

- Labels of  $L_S$  are bound to continuations in  $\rho_S$ . For each label the compile time environment  $y$  gives a corresponding label of  $L_T$ . Suppose, we have label  $N$  in  $L_S$  bound to  $\gamma = \rho_S[N]$ . Assume that according to  $y$  label  $\gamma$  corresponds to  $j = y[N]$  in  $L_T$ . Then we require that  $\gamma$  and  $\rho_{Tj}$  are equivalent with respect to  $\vartheta_G$ .
- Procedure identifiers in  $L_S$  are bound to procedure values in  $\rho_S$ . In the compile time environment a label in  $L_T$  is assigned to each procedure identifier. The continuation bound to this label in  $\rho_T$  has to correspond to the procedure value according to  $\vartheta_P$ .

$$\begin{aligned} \vartheta_V(\rho_S, \gamma, y, \rho_T) \equiv & \forall N \in \text{Num}. \vartheta_G[\rho_S^{\#1}[N], s, \rho_T[y[N]]] \wedge \\ & \text{where } n = \rho_S^{\#5}[I] \\ & \forall I \in Id. \vartheta_P[\rho_S^{\#2}[I], n, m, \rho_T[y[I]]] \\ & \text{where } n = \rho_S^{\#3}[I] \\ & \text{where } s = \Phi \lambda S P Y \\ & \text{where } \{I\} :: [\text{proc}; r_1, \dots, r_m] \end{aligned}$$

#### 5.3.5. Existence of recursive predicates

The definition of the predicates is cyclic or recursive. Therefore we have to investigate whether these relations are at all well defined. We did not have similar problems with recursive functions since our formal language syntactically guarantees that all functions are continuous and for continuous functions least fixed points always exist.

The language used to define predicates  $\vartheta$ , is a significant extension of Scott's logic of computable functions. In particular set constructors are used freely. Therefore there is no static guarantee that the so defined relations (functions into  $T$ ) are continuous. In fact, the predicates defined above are not

even monotonic. For example, we have

$$\begin{aligned} \vartheta_S(\sigma, l, \mu) \equiv & \sigma^{\#2} = \mu^{\#2} \wedge \\ & \sigma^{\#3} = \mu^{\#3} \wedge \\ & \forall a, n. l(a, n) \wedge n > \mu^{\#4} \Rightarrow \vartheta_Y(\sigma a, l, \mu n) \wedge \\ & \forall a, n. l(a, n) \wedge n \leq \mu^{\#4} \Rightarrow \sigma a = \text{unused} \end{aligned}$$

If  $l$  increases  $\vartheta_S$  becomes smaller.

The above situation has been discussed extensively in the literature [MS76, Re74, St77]. It has been shown that under certain conditions recursively defined predicates have a least fixed point even if they are not monotonic. To repeat the theory of "inclusive" predicate is beyond the scope of this thesis. However, Reynolds [Re74] presents a set of easy to check criteria that ensure existence of fixed points. It can be seen that our definitions satisfy his criteria; thus we subsequently assume that the above predicates are well defined.

#### 5.4. Implementation of the code generation

##### 5.4.1. Specifying code generating procedures

The relations introduced above can be used to specify the correctness of our recursive code generating procedures. For example, *Ecode* generating code for expressions, is specified as

```
procedure Ecode(E: assign; f: U; p: Us; y: Uc; rho: Uf; var z: code);
initial z == z0;
exit VrTr: (thetaU(p,f,y,rho) \wedge thetaG(r,PhiExP y, M[z0]rho;r)) =>
thetaG(f[E]zsr, PhiExP y, M[z]rho;r);
```

Note, that this specification assumes that there is some initial sequence of code  $z0$  to which the new code is appended, yielding  $z$ .

One problem with this specification is that it involves quantification. Our solution is to introduce auxiliary parameters to allow instantiation of universally quantified variables as outlined in chapter II. For each expression we know the source continuation which applies. If we also knew the target continuation we could provide the correct instances of these continuations with each call to *Ecode* and thus avoid quantification. The specification would then be

```
procedure Ecode(E: assign; f: U; p: Us; y: Uc; rho: Uf; r: Tr; Gr; var z: code);
entry thetaU(p,f,y,rho) \wedge thetaG(r,PhiExP y, M[z]rho;r);
exit thetaG(f[E]zsr, PhiExP y, M[z]rho;r);
```

occurrence of  $z$  in the entry conditions refers to the initial value of  $z$  in the exit condition refers to the final code.

In the instantiation technique we have to know the correct instantiations of the target continuation. But this in turn requires that we know the expression we are translating. For this reason we decide to generate code backwards! That is, the code generation will produce code in reverse order, last instruction first.

Let us now look at the definition of *AEcode* (*AEcode* corresponds to  $\ell$ , it has specifications identical to *Ecode*). Some typical cases are

```

procedure AEcode(E: asyn;
  zeta: static-environment;
  rho: Tenvironment;
  gamma: S continuation;
  y: compile-environment;
  rhoT: T environment;
  gammat: T continuation;
  var z: code);
entry varthetaU(rho, zeta, y, rhoT) ∧
varthetaC(gamma, Phil(zeta, rho, y), Mscr(z, rhoT, gammat));
exit varthetaC(AEscr(E, zeta, rho, gamma), Phil(zeta, rho, y),
Mscr(z, rhoT, gammat));
var E1, E2, Elst, I, N, O: asyn;
alpha: location;
epsilon: value;
mu, multist: mode;
nu, m: integer;
tau, tau0, tau: type;
begin
  if isnumber(E) then
    begin
      matchnumber(E, N);
      cmkelt(Nscr(N), rhoT, gammat, z);
    end
  else
    if ismid(E) then
      begin
        mu ← ce[E1, zeta];
        tau ← ctyp(mu);
        if isarraytype(tau) then
          begin
            matcharraytype(tau, nu, tau0, tau1);
            indexcode(zeta, rho, gamma, y, rhoT, gammat, z);
            Acode(E2, zeta, mkevalmode(tau0), rho, index(gamma), y, rhoT, gammat, z);
            Ecode(E1, zeta, rho, Aicr(E2, zeta, mkevalmode(tau0), rho,
              index(gamma)), y, rhoT, gammat, z);
          end
        else error
      end
    else if isselect(e) then
      begin
        matchselect(E, E1, I);
        selectcode(I, zeta, rho, gamma, y, rhoT, gammat, z);
        Ecode(E1, zeta, rho, select(I, gamma), y, rhoT, gammat, z);
      end
    else error;
  end;

```

Most of the details of generating code are delegated to auxiliary functions such as *selectcode*, *indexcode*, *cmkelt* and so on. Typical specifications of these functions are

```

procedure cmkelt(epsilon: value; rhoT: Tenvironment;
  gammat: T continuation, var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammat) = Tapply(Mscr(z0, rhoT, gammat), epsilon);
external;
```

```

procedure cmkelnop(O: asyn; rhoT: Tenvironment;
  gammat: T continuation; var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammat) = Tunop(Mscr(z0, rhoT, gammat), O);
external;
```

Let us explain *cmkelt*, which generates an instruction *LIT*  $\epsilon$  and appends

The proof of the verification conditions for *AEcode* is immediate from the definition of  $\hat{\epsilon}$  and some to be defined properties of the predicates  $\vartheta_i$ . Consider the following example for the case where  $E$  is a numeral.

$$\begin{aligned} & (\text{varthetaG}(\gamma, \text{PhiS}(zeta, rho, y), \text{Macr}(z, rhoT, gammaT))) \wedge \\ & \text{isnumber}(e) \wedge \\ & e = \text{makenumber}(n_0) \wedge \\ & \text{Macr}(z_0, rhoT, gammaT) = \text{Tapply}(\text{Macr}(z, rhoT, gammaT), \text{Ncr}(n_0)) \\ \Rightarrow & \text{varthetaaG}(\text{AEscr}(e, zeta, rho, gamma), \text{PhiS}(zeta, rho, y), \\ & \text{Macr}(z_0, rhoT, gammaT))) \end{aligned}$$

For the system's prover this verification condition poses no problem, given a suitable axiomatization of the concepts involved (see appendix 6.). Let us demonstrate the validity of the above condition manually to illustrate the basic schema of proof used to verify the code generation. It is advantageous to rewrite the verification condition in a shorter mathematical notation.

$$\begin{aligned} & \vartheta_G(\gamma, \Phi_{ES} \rho y, M[z] \rho \tau \gamma) \wedge \\ & \vartheta_U(\rho, \varsigma, y, \rho \tau) \wedge \\ & M[z_0] \rho \tau \gamma = \lambda \delta \sigma. M[z] \rho \tau \tau \delta(n \sigma) \\ \Rightarrow & \vartheta_G(\hat{\epsilon}[N] \varsigma \rho \gamma, \Phi_{ES} \rho y, M[z_0] \rho \tau \gamma) \end{aligned}$$

Substituting  $\hat{\gamma}_T$  for  $M[z] \rho \tau \gamma$  and expanding the definition of  $\hat{\epsilon}$  we have to prove (omitting some redundant assumptions)

$$\begin{aligned} & \vartheta_G(\gamma, s, \hat{\gamma}_T) \wedge \\ \Rightarrow & \vartheta_G(\text{apply } \gamma(N[N]), s, \lambda \delta \sigma. \hat{\gamma}_T \delta(N[N] \cdot \sigma)) \end{aligned}$$

This formula is a theorem and can be added to the logical basis of the verification. It can be proven as follows. By definition of  $\vartheta_G(\gamma, s, \hat{\gamma}_T)$  we get

$$\forall X \beta \delta. \vartheta_X(X, \beta, \delta) \Rightarrow \vartheta_G(\gamma X, l, \hat{\gamma}_T \delta)$$

for  $l = \Phi \wedge s \chi \delta \cup \beta^{#3} \cup I_h$ .

Since we have  $\vartheta_V(N[N], l, N[N])$  the definition of  $\vartheta_G$ , yields

$$\vartheta_G(\gamma X, l, \hat{\gamma}_T \delta) \Rightarrow \vartheta_G(\gamma X \chi \delta(N[N], l, \lambda \sigma. \hat{\gamma}_T \delta(N[N] \cdot \sigma)))$$

for any  $X, \beta$ , and  $\delta$  satisfying  $\vartheta_X(X, \beta, \delta)$ . By definition this means that

$$\vartheta_G(\text{apply } \gamma(N[N]), s, \lambda \delta \sigma. \hat{\gamma}_T \delta(N[N] \cdot \sigma))$$

holds.

Similar theorems about properties of  $\vartheta_i$  are required for other cases of *Ecode*. Prove of these theorems is not possible within the verifier as it involves quantification.

Other valuations follow the schema outlined for  $\epsilon$ . However, there are some points that require some further attention.

- What happens to labels? How does the fixed point describing the meaning of labels relate to the produced code?
- How are declarations treated?
- How are procedures translated?

#### 5.4.2. Treatment of labels

The meaning of a sequence of statements is defined by *C $\ell$*  as a least fixed point binding labels to continuations. For the code generating procedure corresponding to *C $\ell$*  we have to prove that the code produced for a list of statements has the correct meaning. To do this we produce a block. The meaning of this block is given as a least fixed point binding labels to target continuations. Thus we have to prove the equivalence of two fixed points, *fix f<sub>S</sub>* and *fix f<sub>T</sub>*. More precisely, corresponding to *C $\ell$*  we use a procedure *Ccode* to produce code for statement sequences. For given  $y$  and  $\varsigma$  we want to prove that *Ccode* satisfied the entry-exit specification:

$$\begin{aligned} & \text{entry } \vartheta_U(\rho, \varsigma, y, \rho \tau) \wedge \vartheta_G(\gamma, \Phi_{ES} \rho y, \gamma \tau); \\ & \text{exit } \vartheta_G(C[\ell[\Gamma]] \rho \gamma, \Phi_{ES} \rho y, M[z] \rho \tau \gamma); \end{aligned}$$

where  $\gamma$  and  $\gamma \tau$  are the correct instances of the source and target continuations.

According to the semantics of *L<sub>S</sub>* and *L<sub>T</sub>* we have

$$\begin{aligned} C[\ell[\Gamma]] \varsigma \rho \gamma &= C[\Gamma] \hat{\varsigma} \hat{\rho} \gamma \\ \text{where } \hat{\varsigma} &= \varsigma[\dots, \text{true}, \dots]/j[\Gamma], \\ \hat{\rho} &= \text{fix}(\lambda \hat{\rho}. \rho [j[\Gamma]] \hat{\varsigma} \hat{\rho} \gamma / j[\Gamma])) \end{aligned}$$

and for  $z = (I_1 \dots I_n)$

$$M[z] \rho \tau \gamma = M^* [I_1 \dots I_n] \rho \gamma$$

where  $\rho_1 := \text{fix}(\lambda \rho. \rho [L[I_1 \dots I_n] \rho \gamma / l[I_1 \dots I_n]])$

Assume  $y$  and  $\varsigma$  are fixed for now and suppose that we can prove the

following property for  $I_1, \dots, I_n$ .

$$\begin{aligned} & \forall \bar{\rho} \bar{\rho}_T \bar{\gamma} \bar{\gamma}_T (\vartheta_U(\bar{\rho}, \bar{\gamma}, \bar{\rho}_T) \wedge \vartheta_G(\bar{\gamma}, \Phi_{\Sigma} \bar{\rho}, \bar{y}, \bar{\gamma}_T) \Rightarrow \\ & \quad \vartheta_G(C[[\Gamma]] \bar{\rho}, \Phi_{\Sigma} \bar{\rho}, y, M[[I_1, \dots, I_n]] \bar{\rho}_T \bar{\gamma}_T) \wedge \\ & \quad \vartheta_G(J[[\Gamma]] \bar{\rho}, \Phi_{\Sigma} \bar{\rho}, y, L[[I_1, \dots, I_n]] \bar{\rho}_T \bar{\gamma}_T)) \end{aligned}$$

where  $\vartheta_G^*$  is true for two lists of continuations if corresponding elements are related by  $\vartheta_G$ . Assuming that the initial environments are related (entry condition), i.e.  $\vartheta_U(\bar{\rho}, \bar{\gamma}, \bar{\rho}_T)$ , then by fixed point induction it is immediate that  $\vartheta_U(\bar{\rho}, \bar{\gamma}, y, \rho_T)$ . Assuming further, that the initial continuations are related (exit condition), i.e.  $\vartheta_G(\bar{\gamma}, \Phi_{\Sigma} \bar{\rho}, y, \gamma_T)$ , the exit condition follows.

In the above arguments we assumed that the predicate  $\vartheta_U$  is admissible. In fact, Reynold's theorem [Re74] mentioned earlier not only guarantees the existence of the recursive predicates  $\vartheta$ , but also proves that these predicates are admissible<sup>3</sup>.

#### 5.4.2.1. Implementation of Ccode

We assume a procedure Ccode, corresponding to C, with the specification

$$\begin{aligned} & \text{procedure } C\text{code}[\Gamma; \text{asym}; \varsigma; U_i; Y; \text{ var } z; \text{code};] \\ & \text{exit } \forall p \forall r \forall u (\rho, \bar{y}, \rho_T) \wedge \vartheta_G(\bar{\rho}, \Phi_{\Sigma} \bar{\rho}, y, \bar{r}) \Rightarrow \\ & \quad \vartheta_G(C[[\Gamma]] \bar{\rho}_T, \Phi_{\Sigma} \bar{\rho}, y, M[[z]] \rho_T \gamma_T) \wedge \vartheta_G(J[[\Gamma]] \bar{\rho}_T, \Phi_{\Sigma} \bar{\rho}, y, L[[z]] \rho_T \gamma_T); \\ & \quad U_i = y J[[\Gamma]]; \end{aligned}$$

Procedure Ccode can then be written as

$$\begin{aligned} & \text{procedure } C\text{code}[\Gamma; \text{asym}; \varsigma; U_i; Y; \text{ var } z; \text{code};] \\ & \text{entry } \vartheta_U(\bar{\rho}, \bar{\gamma}, y, \rho_T) \wedge \\ & \quad \vartheta_G(\bar{\gamma}, \Phi_{\Sigma} \bar{\rho}, y, M[[z]] \rho_T \gamma_T); \\ & \text{exit } \vartheta_G(C[[\Gamma]] \bar{\rho}_T, \Phi_{\Sigma} \bar{\rho}, y, M[[z]] \rho_T \gamma_T); \\ & \text{var } z\text{new}: \text{code}; \\ & \quad \varsigma p U_i; \\ & \quad y\text{new}: Y; \\ & \begin{aligned} & \text{begin} \\ & \quad \varsigma \leftarrow \varsigma \cup \{ \dots, \text{true}, \dots \} / j[[\Gamma]] \} \\ & \quad mkeTLabels(y_j[[\Gamma]], y\text{new}); \\ & \quad C\text{code}[\Gamma, \varsigma, y\text{new}, z\text{new}]; \\ & \quad cmkblockcode(z\text{new}, z, \rho_T, \gamma_T); \\ & \quad \text{end}; \end{aligned} \\ & \text{3.) Inclusive predicates are admissible} \end{aligned}$$

*mkeTLabels* generates a new set of labels for LT corresponding to  $j[[G]]$  and enters this correspondence in  $y$ , resulting in  $y\text{new}$ .

The procedure *cmkblockcode* converts the list of instructions *znew* into a block (*znew*) and appends this block to the front of  $z$ . The parameters  $\rho_T$  and  $\gamma_T$  are virtual and only required in the entry-exit specifications.

procedure *cmkblockcode*(*znew*:code; var  $z$ :code;  $\rho_T$ :U;  $\gamma_T$ :G);  
initial  $z = z0$ ;  
exit  $M[[z]] \rho_T \gamma_T = M[[Tmkeblock(znew)]] \rho_T; M[[z0]] \rho_T \gamma_T$ ;

For the above implementation of Ccode the following verification condition has to be proven.

$$\begin{aligned} & \vartheta_U(\bar{\rho}, \bar{\gamma}, y, \rho_T) \wedge \vartheta_G(\bar{\gamma}, \Phi_{\Sigma} \bar{\rho}, y, M[[z]] \rho_T \gamma_T) \wedge \\ & \quad \text{distinct}(y, j[[G]]) \wedge \\ & \quad \forall \bar{\rho}_T \bar{\gamma}_T (\vartheta_U(\bar{\rho}, \bar{\gamma}, \bar{\rho}_T) \wedge \vartheta_G(\bar{\gamma}, \Phi_{\Sigma} \bar{\rho}, \bar{y}, \bar{\gamma}_T) \Rightarrow \\ & \quad \vartheta_G(C[[G]] \bar{\rho}, \bar{y}, M[[z]] \bar{\rho}_T \bar{\gamma}_T) \wedge \\ & \quad \vartheta_G(J[[G]] \bar{\rho}, \bar{y}, \Phi_{\Sigma} \bar{\rho}, \bar{y}, L[[z]] \bar{\rho}_T \bar{\gamma}_T) \wedge \\ & \quad \ell[[z]] = \bar{y} J[[G]] \\ & \Rightarrow \\ & \quad \vartheta_G(C[\ell][[G]] \varsigma \rho_T \gamma_T \Phi_{\Sigma} \bar{\rho}, y, M[[z]] \rho_T \gamma_T); \\ & \quad \text{where } \varsigma = \{ \dots, \text{true}, \dots \} / j[[G]] \} \\ & \quad \text{where } \bar{y} = y[[\text{label } \ell / j[[G]]]] \} \\ & \quad \text{for new unique labels } \text{label} \end{aligned}$$

We let  $\ell_T = M[[z]] \rho_T \gamma_T$ . If we apply the definitions of C $\ell$  and M to the conclusion of the VC we get

$$\begin{aligned} & \vartheta_G(C[[G]] \bar{\rho}, \bar{y}, M[[z]] \rho_T \gamma_T) \\ & \text{where } \bar{\rho} = \text{fix } (\bar{\rho}, \rho)[J[[G]] \bar{\rho}, \gamma_T J[[G]]] \\ & \text{where } \bar{\rho}_T = \text{fix } (\bar{\rho}_T, \rho_T)[L[[z]] \bar{\rho}_T \gamma_T J[[G]]] \} \end{aligned}$$

Since locations bound in  $\varsigma$ ,  $\rho$ , and  $y$  are identical to those in  $\bar{\varsigma}$ ,  $\bar{\rho}$ , and  $\bar{y}$  we have

$$\Phi_{\Sigma} \bar{\rho} \bar{y} = \Phi_{\Sigma} \bar{\varsigma} \bar{\rho} \bar{y}$$

Now, the proof of this verification condition follows by fixed point induction as outlined above.

#### 5.4.2.2. Weak V-introduction

The remaining problem is to prove the universally quantified exit condition of Ccode.

$\forall \rho \rho \tau \tau \vartheta_U(\rho, s, y, \rho \tau) \wedge$   
 $\vartheta_G(C[[G]]\rho \tau, \Phi_E \rho y, M[[z]]\rho \tau \tau) \Rightarrow$   
 $\vartheta_G(C[[G]]\rho \tau, \Phi_E \rho y, L[[z]]\rho \tau \tau)$

The solution is weak  $\forall$ -introduction defined in chapter II. We prove that the body of Ccode satisfies the following quantifier free specifications.

procedure Ccode(G:gamma; {U; rho; U\_r; var z:code};  
entry  $\vartheta_U(\rho, s, y, \rho \tau) \wedge \vartheta_G(\rho, \Phi_E \rho y, \tau)$ ;  
exit  $\vartheta_G(C[[G]]\rho \tau, \Phi_E \rho y, M[[z]]\rho \tau \tau) \wedge \vartheta_G(J[[G]]\rho \tau, \Phi_E \rho y, L[[z]]\rho \tau \tau)$ ;

The rule of weak  $\forall$ -introduction allow to use the specifications

entry true;  
exit  $\forall \rho \rho \tau \tau \vartheta_U(\rho, s, y, \rho \tau) \wedge \vartheta_G(\rho, \Phi_E \rho y, \tau) \Rightarrow$   
 $\vartheta_G(C[[G]]\rho \tau, \Phi_E \rho y, M[[z]]\rho \tau \tau) \wedge \vartheta_G(J[[G]]\rho \tau, \Phi_E \rho y, L[[z]]\rho \tau \tau)$ ;

for calls to Ccode. The universally quantified exit condition can be abbreviated as

$\forall allU GG(G, s, y, z) \equiv$   
 $\vartheta_U(\rho, s, y, \rho \tau) \wedge \vartheta_G(\rho, \Phi_E \rho y, \tau) \Rightarrow$   
 $\vartheta_G(C[[G]]\rho \tau, \Phi_E \rho y, M[[z]]\rho \tau \tau) \wedge \vartheta_G(J[[G]]\rho \tau, \Phi_E \rho y, L[[z]]\rho \tau \tau)$

in our assertion language.

We have shown above that by fixed point induction it is possible to prove

$$\begin{aligned}
&\vartheta_U(\rho, s, y, \rho \tau) \wedge \vartheta_G(\rho, \Phi_E \rho y, M[[z]]\rho \tau \tau) \wedge \\
&\text{distinct}(y[[G]]) \wedge \\
&\text{forallUG}(G, \dot{y}, z) \wedge \\
&\dot{z}[z] = y[j[G]] \\
\Rightarrow &\vartheta_G(C[[G]]\rho \tau, \Phi_E \rho y, M[[z]]\rho \tau \tau) \\
\text{where } &\dot{s}[\dots, \text{true}, \dots, j][G] \\
\text{where } &y[[\text{abs}/j[G]]]
\end{aligned}$$

for new unique labels  $\text{abs}$

This theorem can now be added to our logical basis as a rule

```

infer varthetaG(C as if(G, zeta, rho, gamma), s, Macr[7mkeblock(z), rhoT, gammaT])
from forallUG(G, rho, true, j, rhoT)
varthetaU(rho, zeta, y, rhoT) \wedge
s = PhiS(zeta, rho, y) \wedge
disjointLabels(y, j)(G), y, \wedge

```

$V \text{arthetaG}(\text{gamma}, s, \text{gammaT})$ ;

#### 5.4.3. Declarations

There are code generating procedures corresponding to valuations that deal with declarations (e.g.  $V, Q$ ). But these procedures will not generate code, rather they will change the compile time environment by allocating new storage.

For example we have

$$V_U[I:T]\varsigma\rho = \text{if } [\text{var}:r]::[I] \text{ then } \rho[\alpha/I]$$

where  $(\alpha, \rho) = \text{new}\ r\ \rho$ .

The corresponding procedure *Allocate* will allocate memory for  $I$  in  $L.T$ .

```

procedure Allocate(D:asym; {U; rho; U_r}; var y:Y; rho; U_r);
entry  $\vartheta_U(\rho, s, y, \rho \tau)$ ;
exit  $\vartheta_U(V, [I:T]\rho, s, y, \rho \tau)$ ;
begin
  if [ $\text{var}:r$ ]:{[]/I} then
    begin
      n  $\leftarrow$  size(r);
      y  $\leftarrow$   $y^{n+1}/I, y^{n+2}, y^{n+3}, y^{n+4} + n$ ;
    end
  end;

```

A proof of this procedure is immediate, considering that no continuation or procedure value in  $\rho$  or  $\rho \tau$  have been changed.

Meaning functions that describe initialization of variables are similar to other executable statements; their implementation follows analogously.

#### 5.4.4. Procedures and functions

We mentioned above that there are no procedure and function values in  $L.T$ . A procedure in  $L.S$  corresponds to a label in  $L.T$ , the start address of the code for this procedure. The relation  $\vartheta_U$  guarantees that for any procedure value  $\pi$  and the associated continuation  $\gamma$  of the start address of the procedure the following holds:

$$\forall \gamma_S \gamma_T \vartheta_G(\gamma_S, \gamma_T) \Rightarrow \vartheta_G(\pi \gamma_S, \dot{\gamma})$$

where  $\dot{\gamma} = \lambda(\sigma, \gamma)(\text{up } \delta \gamma n(\text{length } \sigma - m)k)\sigma$

Here,  $n$  is the lexical level of the definition of the procedure,  $m$  is the number of parameters, and  $k$  is the number of cells allocated in the calling environment.

Since the semantics of the *CALL* instruction is given by

$$\mathcal{M}[\text{CALL } l \ n \ m \ k] \rho \gamma = \lambda \delta \sigma. \rho[\{\text{up } \delta \ n(\text{length } \sigma - m)k\} \sigma]$$

the correctness of code for a procedure call is immediate. The necessary parameters to *CALL* are readily available at call time.  $k$  is given as  $y^{\#4}$ .

It remains to be shown that correct code is generated for a procedure declaration. We have the definition

$\mathcal{F}[\text{procedure } I(\Pi_1, \dots, \Pi_n); \Theta] \rho = (\pi)$

where  $\pi = \lambda \gamma. \text{enter}(n + 1); P[\Pi_1, \dots, \Pi_n]_{\varsigma_1} \rho_i$ ;

$B[\Theta]_{\varsigma_2} \rho_i; \text{exit } \gamma$

where  $\varsigma_2 = \varsigma_1[\text{proc } \mu_1, \dots, \mu_n / I]$

where  $(\mu_1, \dots, \mu_n), \varsigma_1 = P[\Pi_1, \dots, \Pi_n]_{\varsigma_1}$

where  $\rho_1 = Q[\Pi_1, \dots, \Pi_n]_{\varsigma_1}(\text{next } \rho)$

where  $n = \text{level } \rho$

Code for such a declaration is generated by a procedure *Fcode* with the following specifications.

*procedure*  $Fcode(P, \text{asym}, \varsigma, U, \rho, U_r, G_r, \rho_r, U_r, \text{var } z: \text{code};$   
 $\text{exit } \vartheta_P(\mathcal{F}[P] \rho, \text{level } \rho, \text{params}(P, \varsigma), M_{\varsigma}[\rho_r]_{\rho_r})$

where *params* returns the number of parameters of  $P$  in the environment  $\varsigma$ . *Fcode* generates the following code.

- *Qcode* is called which changes the compile time environment to contains bindings of parameters to target locations. Let  $\hat{\varsigma}$  be this new compile time environment.
- *Pcode* generates code  $z_P$  that initializes parameters with the values of actual parameters. These values are all on top of the stack (in reverse order).
- *Bcode* produces code  $z_B$  for the procedure body.
- The *EXIT* instruction causes return from the procedure.

According to the semantics of *L7* we have  $\mathcal{M}[\text{EXIT}] \rho_r \gamma r = \lambda \delta. \delta^{\#5} \delta^{\#3}$ . Thus, for the final code  $z$  produced by *Fcode* we have

$$\mathcal{M}[z] \rho_r \gamma r = \mathcal{M}[z_P; z_B] \rho_r; \lambda \delta. \delta^{\#5} \delta^{\#3}.$$

By the source semantics we have

$$\mathcal{F}[P] \varsigma \rho = \lambda \gamma. \text{enter}(n + 1); P[\Pi]_{\varsigma} \rho_i; B[\Theta]_{\varsigma_2} \rho_i; \text{exit } \gamma$$

where  $\varsigma_1$ ,  $\varsigma_2$  and  $\rho_1$  are defined as in the definition above;  $n = \text{level } \rho$ .

Expanding the exit condition of *Fcode*, we have to prove that

$$\begin{aligned} &\forall s, \gamma s, \hat{\gamma}_T. \vartheta(\gamma s, \hat{\gamma}_T) \Rightarrow \vartheta_G(\bar{\gamma}_S, s, \bar{\gamma}_T) \\ &\text{where } \bar{\gamma}_S = \text{enter}(n + 1); P[\Pi]_{\varsigma_1} \rho_i; B[\Theta]_{\varsigma_2} \rho_i; \text{exit } \gamma S \\ &\text{where } \bar{\gamma}_T = \lambda \delta \sigma. \mathcal{M}[z_P; z_B] \rho_T(\lambda \delta. \delta^{\#5} \delta^{\#3})(\text{up } \delta \hat{\gamma}_T \ n(\text{length } \sigma - m)k) \sigma \end{aligned}$$

where  $k = \text{used}(s)$ .

To prove this verification condition we use the fact that

$$\mathcal{M}[z] \rho_T \gamma r \delta = g[z] \rho_T(\gamma r \delta) \delta$$

for some function  $g$ . This formula is not valid for arbitrary pieces of code. However, it holds for all code sequences  $z$  produced by our code generating functions.

Now, consider  $\mathcal{M}[z] \rho_T(\lambda \delta. \delta^{\#5} \delta^{\#3}) \delta$ . By the above lemma we have

$$\mathcal{M}[z] \rho_T(\lambda \delta. \delta^{\#5} \delta^{\#3}) \delta = g[z] \rho_T((\lambda \delta. \delta^{\#5} \delta^{\#3}) \delta) \delta$$

$$g[z] \rho_T((\lambda \delta. \delta^{\#5} \delta^{\#3}) \delta) \delta = \mathcal{M}[z] \rho_T(\lambda \delta. \delta^{\#5} \delta^{\#3}) \delta$$

Applied to the verification condition for *Fcode* we can rewrite  $\bar{\gamma}_T$  as

$$\bar{\gamma}_T = \lambda \delta \sigma. \mathcal{M}[z_P; z_B] \rho_T; \lambda \delta. \delta \hat{\gamma}_T \ n(\text{length } \sigma - m)k \sigma$$

Thus, we have shown, that the procedure returns to the right point.

The proof of the verification condition now follows from the correctness of the components  $z_P$  and  $z_B$ . With  $\vartheta_G(\gamma_S, s, \bar{\gamma}_T)$  we get  $\vartheta_G(\gamma_S, \delta, \lambda \delta. \delta^{\#5} \delta^{\#3})$ . From this and the correctness of  $z_B$  we get

$$\vartheta_G(B[\Theta] \varsigma \rho; \text{exit } \gamma, \mathcal{M}[z_B] \rho_T(\lambda \delta. \delta^{\#5} \delta^{\#3}))$$

and so on.

#### 5.4.5. Blocks

The meaning of a block in *L7* is defined by a least fixed point similar to *C7*. The main difference is that not only labels are bound in the environment by also procedures and functions are bound to their procedure and function values. The implementation strategy for blocks is the same as that for command lists. Weak  $\forall$ -introduction is used in the very same way as for *Ccode*.

#### 5.4.6. Refinement

We have described the overall strategy as well as all of the non obvious parts of the code generation. The remaining refinement steps follow along the

same line as for static semantics. We use the same representations for the abstract syntax, type, modes, and environments.

We do not refine the representation of code sequences any further here.

Using the techniques outlined earlier this step should be simple given a concrete representation for the machine code.

Some typical code generation procedures are given in appendix 6.

### 1. Summary

We have presented the theory, specification, implementation, and correctness proof of a compiler. The source and target languages  $L_S$  and  $L_T$  are realistic and useful languages rather than toys.

We have shown that a correct compiler can be systematically developed from a formal definition of source and target languages. The techniques employed in this thesis are very general and will also be applicable to other large software systems.

We have given a formal denotational definition of the source language  $L_S$  on a low level of abstraction. This definition captures and formalizes standard compiling techniques. Future compilers will profit from these results whether or not they are formally verified.

In the course of the proofs we had to tackle various technical problems, such as the treatment of pointer operations, quantification in a quantifier free assertion language and the treatment of fixed points. Though apparently minor, these problems are of general importance and will most certainly arise in totally different applications.

We hope that this work is a convincing argument, that today's verification techniques are sufficiently powerful to be used in real life software development. At the same time this work reveals some of the weak points of the technology and points to future areas of research; we will discuss some possible improvements in the following sections.

Certainly there is no guarantee that our compiler will never fail. We have discussed some of the possible sources of errors earlier. Still, software developed from formal specifications and formally verified increases our confidence in its correctness to a point not achievable with conventional testing techniques alone. As the development of techniques for program verification will continue future proofs will leave less and less margin for errors.

Our compiler is not yet a "software product"; several features have to be added to make it a useful compiler. A suitable error handling and recovery mechanism as well as output of the produced code in suitable format have to be added. As we pointed out earlier, these additions can be made without invalidating the program proofs given. Furthermore, it is desirable to include a minimum of code optimization in the compiler. We will discuss some of the possibilities for this as well as other useful extensions in the following section. Numerous programs are written today for which correctness is imperative. We believe that in many cases verification is economically feasible and that these programs would greatly benefit from verification.

## 2. Extensions

It is desirable to be able to verify compilers that involve code optimizations. This is particularly important since a large number of the errors of present compiler are in the optimization part of the compiler.

Also, one may wish to apply our techniques to languages with features not present in *L.S.* Both points are discussed in this section.

### **2.1. Optimization**

Except for compile time evaluation of constant expressions our compiler does not perform optimizations of any kind. Let us briefly consider how possible optimizations would affect our proofs.

Let us classify optimizations in those which require data flow analysis and those which do not. In both cases we can distinguish optimizations performed on the source and the target level.

Optimizations that require data flow analysis are very hard to include in our compiler. They require extensive proofs that certain manipulations of the source or target text do not alter the meaning of a program. The correctness of these transformations may depend on an arbitrarily large context. The necessary theorems cannot be proven (or even formulated) in our assertion language. A reasonable approach would be to include an additional "optimization" step in the compiling process and develop a suitable logical basis for it. But it should be expected that a sophisticated optimization alone approaches or exceeds the complexity of the complete compiler presented here.

Optimizations that do not require flow analysis present a much brighter picture. For example, given the sequence of assignments  $x \leftarrow a; z \leftarrow b$  where

$a$  and  $b$  are free of side effects (say simple variables) then we can omit the first assignment without changing the meaning of the program. Sequences of this nature can be detected easily and the proof that certain optimizations preserve the meaning of the program is easy. Unfortunately, redundant operations are not very common in the source language.

By far the easiest optimization to include in our compiler is a technique called "peephole optimization" introduced by McKeeman [Mc65]. The idea is to find redundant operations in the produced target code or to find sequences of code that can be replaced by shorter code. The idea is very simple. One moves a window across the final code and matches the instructions in the window with a given set of patterns for which a correctness preserving transformation is known.

The technique will be very successful since our code generation is very primitive. Code for components of the source program is produced irrespective of the context. For example, code sequences of the form  $HOP\ n, LABSET\ n$  are very common. Also, our instruction set was chosen to be minimal. Most machines have instructions that are not strictly necessary but speed up certain frequent operations. *L.T* could have an instruction *INCR* which increments the top of the stack. Then the optimizer could search for the pattern *LIT 1; ADD* and replace it by *INCR*.

The verification of a peephole optimizer is fairly straightforward. One can prove the validity of a set of transformations of the form  $pattern_1 \rightarrow pattern_2$  and systematically apply them.

### **2.2. Register machines**

Generating code for a register machine is substantially more complicated than producing code for a stack machine. Still, a sensible approach would be to first produce intermediate code for a stack machine and in a later translation step to convert it into code for a register machine. This latter translation can be isolated from the rest of the compiler in which case a correctness proof becomes manageable. Sophisticated code for a register machine will require optimizations (e.g. register allocation) and the remarks of the previous sections apply.

### **2.3. New language features**

Let us first discuss the addition of data structures of Pascal omitted in *L.S.* Types *real* and *char* do not cause any problems at all. They merely require quite changes in the definition of *Ty*, the domain of types. Reals require some

attention in the treatment of coercion operations but a solution is straightforward. Similarly, a simple minded implementation of variant records is easy. To model the (undesirable) semantics of Pascal we simply define the location corresponding to records to have several identifiers mapping to the same location. Even simpler is the implementation that assigns disjoint storage to all variants.

Packed data structures pose an interesting problem. Specifying an object to be packed does not change the semantics in any way. The *packed* attribute is of pragmatic rather than semantic nature. Since a denotational definition is purely extensional an implementation ignoring the packed attribute is perfectly valid. Conversely, a compiler could pack all data objects and would still be correct. It is unclear how attributes like *packed* can be captured by a formal definition or whether they are desirable at all.

We omitted for loops, case and with statements in LS. Including any of these poses no problem at all. The case statement and the for loop could both be implemented merely by changing the tree transformations to convert it into nested conditionals and while loops respectively. A simple implementation of the with statement would require to open a new scope in which identifiers denoting record components are declared as variables. In the target language storage would be set up for these variables and their values (addresses) determined on entry to the while statement. In either case a correctness proof is immediate.

Another restriction in LS is that mutual recursion is disallowed. As we mentioned earlier, the code generation translates mutually recursive procedures and functions correctly. The restriction lies in the static semantics part. An extension to cover mutual recursion could be implemented in either of two ways. First, we could have forward declarations as in Pascal. This merely extends the syntax and the tree transformation part of the definition. Changes to the static semantic definition are minor. Alternatively, mutually recursive procedures could be treated similarly to recursive types, i.e. we could allow a procedure or function identifier to occur before its declaration. This requires that in the static semantic definition a fixed point is used to describe the meaning of a set of declarations.

#### 2.4. A stronger correctness statement

We proved that whenever the compiler terminates without printing an error message then the produced code is correct. We argued earlier that is a useful statement since it will never deceive the user of the compiler to believe

that his program is correctly translated when it is not. The reader may ask how difficult it is to prove a stronger statement.

Let us first consider termination. There are well known methods to prove termination in a weak programming logic by adding counters to the program which decrease with every loop and recursion and for which one proves that they always remain positive [LS75, MP73]. It poses no particular problem to include such counters. However, additional program documentation is necessary to prove that they remain positive. For example, consider the implementation of the static semantics. To prove that the recursive procedures terminate requires to prove that the list structure representing the abstract syntax is free of cycles. This is not trivial and requires additional proofs of the tree transformation program to show that only cycle free list structures are produced.

Another possible extension is to prove that whenever the compiler prints an error message, then the input program is wrong. This is true for most error messages. For example, the static semantics defines exactly when an error is called for. Similarly, from the LR-theory it is easily provable that an error message in the parser is only printed if the input program is in error. But at several places in the compiler error messages are printed in situations which are never expected to arise. For example, in the static semantics of expressions we consider all cases of expressions and if none of these cases obtains, an error message is printed. If the abstract syntax tree is build correctly, this situation will never arise and this error message will never occur. However, to prove that the abstract syntax is well formed is not immediate.

In some situations it will be impossible to guarantee, that no errors occur. For example, the stacks used in the parsing algorithm all have finite length. But for any fixed size of these stacks there will be an input program requiring a greater size.

Providing error messages for situations which should never occur or which indicate overflow conditions is standard practice in programming. Our approach shows that these "redundancies" have their place in the framework of verification as well.

#### 3. Future research

This thesis raises several open questions and points to research areas to improve verification. In this section we discuss some of the important issues; but clearly, the list is open ended.

### 3.1. Structuring a compiler

Did we choose the best formal definition? Is our semantic definition the best suited for a compiler proof? Should the structure of the compiler be different? A definite answer to these questions can only be given if we have other verified compilers and can compare different approaches. But some minor points should be noted here.

In the static semantics we computed the modes of all expressions to check the program for validity. This information has been discarded afterwards. Later, the code generation also requires information about modes of expressions; our compiler recomputes modes if necessary. But clearly, this is redundant as one could store the mode information in the abstract syntax tree. The main reason why we choose not to do so is a technical one. Changing the abstract syntax tree is an update operation on a pointer structure which has no obvious counterpart in the formal definition. Consequently it is not immediately clear how the correctness of such update operations can be stated. Alternatively, static and dynamic semantics could be combined in one program. This also eliminates redundant computation of modes. In return the structure of the resulting compiler is less clean.

The idea of combining different parts into one can be carried a step further. The scanner and the parser could be combined into one program thus eliminating the representation of the program as sequence of tokens. A natural way towards this implementation are coroutines which unfortunately are not available in our implementation language. Any two programs that communicate by writing and reading one file can be combined using coroutines without any change in the program proof.

The size of programs that can be translated by our compiler is severely limited by the fact that the complete abstract syntax tree is stored in memory. But this is not strictly necessary. The semantics of *LS* is such that semantic checking and code generation can proceed from left to right requiring only a limited context. A suitable reorganisation of our program seems possible without significant changes in the formal theory. Again, here is a tradeoff between efficiency and structural clarity.

### 3.2. Improvements of verification systems

The research presented here may well be the first real test of any automatic verification system. Several weak points that call for improvement have been revealed. Clearly, the assertion language used in a verifier is a compromise between

expressive power and efficiency of the system's theorem prover. For example, allowing quantification in assertions will make theorem proving much harder and may render the whole system useless.

Still it seems desirable to allow for limited quantification. Can a proof rule corresponding to our weak  $\forall$ -introduction be included. Even if the theorem prover does not deal with quantification, can quantified formulas be allowed syntactically. For example, one might want to write  $\forall z.P(x,y)$  and have the theorem prover treat the formula as a predicate symbol with one free variable  $y$ . Of course, this would not change anything in the capabilities of the system except for the readability of the documentation.

Another desirable feature is a typed assertion and rule language. Many proofs given here were not forthcoming because of a misspelling of a predicate or function symbol. Simple minded type checking could detect many of these errors.

The verifier should provide for virtual code. Not merely virtual variables but rather virtual types, virtual parameters, virtual procedures and so on. Such a concept is imperative if the verified program is ever to be compiled. Also, it greatly enhances the readability of programs and clarifies which parts of the code require further refinement.

The theorem proving part of the Stanford verifier performed extremely well in most cases. In particular the built-in complete decision procedures for frequently occurring theories proved to be a valuable asset. The main problems are in the area of rules supplied to the prover.

The semantic contents of a rule should be stated independently of a particular heuristics of how to use this rule. Currently, there are several different ways of stating the same fact. To get the prover to simplify a verification condition efficiently, rules have to be stated in the "right" format. This requires some experimenting and frequent rewriting of rules, a possible source of errors. In some cases it is desirable to manually intervene in a proof and give hints to the system. Or, one may want to find out why exactly a verification condition does not simplify to true. But manual proofs should be the exception rather than the rule. Giving manual proof guidance to the system in all cases would make this research impossible. Rather what is needed is a general concepts of a "proof schema", a way of telling the system how to prove a verification condition. Consider the static semantics for expressions. The program consists of several branches each of which checks a particular syntactic class of expressions such as identifiers, unary operators, binary operators and so on. The structure of the proof for each of these cases is very similar. Therefore, given the proof of one case, all the other cases follow by analogy.

The concept of a proof schema seems very useful here. We have shown that refinement of a program leads to structurally similar verification conditions. An intelligent verification system might remember previous proofs, thus greatly reducing the amount of work necessary to verify a refined program.

### 3.3. Better verification techniques

Weak  $\forall$ -introduction proved a useful concept in this research. Is it generally applicable? Is it a special concepts or can it be generalised? For example, is there weak  $\exists$ -introduction and so on?

We used operationalization of fixed points to compute recursive types. Is this just a gimmick to circumvent limitations of the verifiers assertion language or is it a generally applicable technique. What is the most general situation in which operationalization can be employed?

Even more interesting is the opposite question. Clearly, any update operation on pointers can be expressed as a least fixed point. Is this a practical way to reason about pointer operations? It certainly is, if pointers are used to represent recursive domains and if the pointer operations have a natural counterpart (e.g. fixed point) in the underlying theory. But can it also be used to prove a tree traversal algorithm which manipulates link fields to avoid stacking operations?

Alternatively, it may be possible to develop high level concepts to reason about pointers similar to those proposed by Reynolds [Re79].

### 3.4. Program development systems

The ideas put forth so far all assume that the current paradigm of program verification is the best possible. Maybe, it is not. What would the ideal verification system for our compiler proof look like? Ideally, of course, we want a program synthesis system, but let us be a little more realistic.

One of the main problems in carrying out the development of the compiler and its proof was consistency. We had to deal with an extremely large formal definition, theorems derived from this definition, machine readable versions of specifications and theorems, several versions (of different refinement level) of the program to be written, and programs communication with a given program. A system that keeps track of all these formal objects and insures notational consistency seems easy to conceive and would be of great assistance.

But, such a system could do even more. We argued previously, that many proofs of cases of our program follow by analogy. But the implementation of

different clauses of the definition follow by analogy also. Thus, what is needed is an intelligent editor. For instance, the user should be able to explain to such a system, how a recursively defined function is to be transformed into an efficient program. Part of such a transformation is the specification of implementation details, for instance, how a particular abstract object is to be represented in the program. An even more advanced system might know about efficient default implementation of certain objects (automatic data structure selection).

References

- Abbreviations:**
- AIM — Artificial Intelligence Memo, Stanford University
  - TCS — Theoretical Computer Science
  - LNCS — Lecture Notes in Computer Science, Springer Verlag
  - CACM — Communications of the ACM
  - JACM — Journal of the ACM
- [AU72] Aho, A. V., Ullman J. D.: *Theory of Parsing, Translation, and Compiling, vol 1, 2*; Prentice Hall, Englewood Cliffs, 1972
- [AJ74] Aho, A. V., Johnson S. C.: *LR Parsing*; Comp. Surveys, Vol. 6, No. 2 (1974)
- [AA74] Aiello, L., Aiello, M., Weyhrauch, R.: *The Semantics of Pascal in LCF*; AIM-221, '74
- [AA77] Aiello, L., Aiello, M., Weyhrauch R.W.: *Pascal in LCF: Semantics and examples of proof*; TCS 5, 1977, 135-177.
- [AE73] Anderson, T., Eye, J., Hornig, J.J.: *Efficient LR(1) Parsers*; Acta Informatica 2 (1973)
- [BL70] Birkhoff, G., Lipson, J. D.: *Heterogeneous Algebras*; J. Comb. Theory 8, 115-133 , 1970
- [Bo76] Bochman, G.: *Semantic evaluation from left to right*; CACM 19/2 (1976)
- [BM77] Boyer, R.S., Moore, J.S.: *A computer proof of the correctness of a simple optimizing compiler for expressions*; SRI Tech. Report 5, 1977
- [Bj77] Bjørner, D.: *Formal development of interpreters and compilers*; Report ID673, Technical University of Denmark, Lyngby 1977
- [BJ78] Bjørner, D., Jones, C.B.(ed.): *The Vienna Development Method: The Meta Language*; LNCS 61 (1978)
- [Bu69] Burstall, R. M.: *Proving Properties of Programs by Structural Induction*; Computer J. 12/1, 1969
- [BL69] Burstall, R. M., Landin, P. J.: *Programs and their Proofs: An Algebraic Approach*; Machine Intelligence 4, 1969
- [BC71] Burroughs Corp.: *Burroughs B6700 Handbook, Vol 1, Hardware*; Form No. 5000276, 1971
- [BC72] Burroughs Corp.: *Burroughs B6700 Information Processing Systems, Reference Manual*; Form No. 1058633, 1972
- [BC73] Burroughs Corp.: *System Miscellanea*; Form No. 5000367, 1973
- [Ca76] Cartwright, R.: *Practical Formal Semantic Definition and Verification Systems*; AIM-296, December 1976
- [CM79] Cartwright, R., McCarthy, J.: *First order programming logic*; Sixth annual ACM Symposium on Principle of Programming languages, San Antonio, 1979
- [Ch76] Chirica L.: *Contributions to compiler correctness*; PhD Diss. UCLA 1976
- [CM75] Chirica, C.M., Martin, D.F.: *An approach to compiler correctness*; Int'l Conf. on reliable software, Los Angeles, 1975
- [Ch51] Church, A.: *The Calculi of Lambda-Conversion*; Annals of Mathematical Studies 6, Princeton University Press, Princeton (1951)
- [Cl79] Clarke E. M.: *Programming language constructs for which it is impossible to obtain good Hoare-like axiom systems*; JACM, 26/1, 1979
- [CH79] Cohen, R., Harry, E.: *Automatic generation of near-optimal linear-time translators for non circular attribute grammars*; Sixth annual ACM Symposium on Principle of Programming languages, San Antonio, 1979
- [Co79a] Cohn, A.: *High level proof in LCF*; Proceedings of the Fifth Symposium on Automated Deduction, 1979
- [Co79b] Cohn, A.: *Machine assisted proofs of recursion implementation*; PhD Thesis, University of Edinburgh, 1979
- [Co65] Cohn, P.M.: *Universal Algebra*; Harper & Row, 1965
- [Co77] Cook, S.: *Soundness and completeness of an axiom system for program verification*; 9th Symp. on Theory of Computing, Boulder, 1977

- [DL78] Darringer, J.A., Laventhal, M.S.: *A Study of the Use of Abstractions in Program Specification and Verification*; IBM Research Report RC1784, 1978
- [DL79] DeMillo, R.A., Lipton, R.J., Perllis, A.J.: *Social processes of proofs of theorems and programs*; CACM 22/5, 1979
- [De78] Deransart, P.: *Proof and Synthesis of Semantic Attributes in Compiler Definition*; IRIA, Report 333, Dec. 1978
- [De71] DeRemer, F.L.: *Simple LR(k) grammar*; CACM 14 (1971)
- [De74] DeRemer, F.L.: *Transformational Grammars for languages and Compilers*; Technical Report 50, Univ. of Newcastle upon Tyne, 1973
- [Di76] Djikstra, E.W.: *A discipline of Programming*; Prentice-Hall, Englewood Cliffs, 1976
- [Do76] Donahue, J. E.: *Complementary Definitions of Programming language Semantics*; LNCS 42, Springer, 1976
- [Do77] Donahue, J.: *On the semantics of "Data Type"*; TR 77-311, CS Dept. Cornell University (1977)
- [DG79] Donzadu, V., Kahn, G., Krieg-Brückner, B.: *Formal definition of ADA*; Preliminary Draft, CII Honeywell-Bull, Oct. 1979
- [En72] Engelriet, J.: *A note on infinite trees; Information Processing Letters* 1, 1972
- [Fl67] Floyd, R. W.: *Assigning Meanings to Programs*; Proceedings of Symp. in Applied Mathematics 19 (1967)
- [Go78] Goguen, J.: *Some Ideas in Algebraic Semantics*; Proceedings of the third IBM Symposium on mathematical foundations of computer science, 1978
- [Gr77] Goguen, J.A., Thatcher, J.W., Wagner, E.G., Wright, J.B.: *Initial algebra semantics and continuous algebras*; JACM 24/1, 1977, 68-95.
- [GH74] Goos, G., Hartmannis, J.: *Compiler Construction, an advanced course*; LNCS 21 (1974)
- [Go75] Gordon, M.: *Operational Reasoning and denotational Semantics*; AIM-264, 1975
- [GM75] Gordon, M., Milner, R., Wadsworth, C.: *Edinburgh LCF*; Internal report CSR-11-77, University of Edinburgh
- [Go79] Gordon, M.J.C.: *The denotational description of programming languages, an introduction*; Springer Verlag, New York, Heidelberg, Berlin, 1979
- [Gr71] Greif, L. A., Meyer: *Specifying Programming language Semantics*; Sixth annual ACM Symposium on Principle of Programming Languages, San Antonio, 1979
- [Gr75] Gries, D.: *Compiler Construction for Digital Computers*; John Wiley, New York, 1971
- [Gu75] Guttag, J. V.: *The Specification and Application to Programming of Abstract Data Types*; Technical Report CSRG- 59, University of Toronto, 1975
- [Gu76] Guttag, J. V., Horowitz, F., Musser, D.R.: *Abstract Data Types and Software Validation*; USC-ISI Technical Report (1976)
- [Ha73] Habermann, A.N.: *Critical Comments on the Programming Language Pascal*; Acta Informatica 3, pp47-57 (1973)
- [He75] von Henke, F. W.: *On the Representation of Data Structures in LCF with Applications to Program Generation*; AIM-267, September 1975
- [HL74] von Henke F. W., Luckham D. C.: *Automatic Program Verification III: A Methodology for Verifying Programs*; AIM-256, 1974
- [Ho69] Hoare, C. A. R.: *An Axiomatic Basis of Computer Programming*; CACM 12, Oct, pp 576-580 (1969)
- [Ho72] Hoare, C. A. R.: *Proofs of Correctness of Data Representation*; Acta Informatica 1/1 (1972)
- [HL74] Hoare, C. A. R., Lauer, P.F.: *Consistent and Complimentary Formal Theories of the Semantics of Programming languages*; Acta Informatica 3, pp135-154, (1974)

- [HW73] Hoare, C. A. R., Wirth, N.: *An Axiomatic Definition of the Programming language Pascal*; Acta Informatica, 2 (1973), pp.335-355
- [Ic79] Ichbiah, J. D. et al: *Preliminary ADA reference manual*; Sigplan Notices 14/6, 1979
- [Il75] Igarashi, S., London, R. L., Luckham, D. C.: *Automatic Program Verification I: Logical Basis and Its Implementation*; Acta Informatica, Vol. 4, pp 145-182 (1975)
- [JW76] Jensen, K., Wirth, N.: *PASCAL, User Manual and Report*; Springer, New York, Heidelberg, Berlin, 1976
- [Jo79] Jones C.B.: *Constructing a theory of a data structure as an aid to program development*; Acta Informatica 11, 1979, 119-137.
- [Ka76] Kaplan, D. M.: *Correctness of a Compiler for Algol-like Programs*; AIM - 48, Stanford University, 1967
- [Ka76a] Kastens, U.: *Systematische Analyse semantischer Abhängigkeiten*; Informatik-Fachberichte 1, Springer Berlin Heidelberg New York, 1976
- [Ka76b] Kastens, U.: *Ein überutzer-erzeugendes System auf der Basis attributierter Grammatiken*; Dissertation, Universität Karlsruhe, 1976
- [Kl67] Kleene, S.C.: *Mathematical Logic*; John Wiley, 1967
- [Kn65] Knuth, D.E.: *On the translation of languages from left to right*; Information and Control 8/6, 1965
- [Kn68] Knuth, D. E.: *Semantics of Context-Free languages*; Math. Systems Theory, vol. 2, pp. 127-145, 1968
- [Kn78] Knuth, D. E.: *Tau Epsilon Chi, A System for Technical Text*; AIM-317, September 1978
- [Kr74] Kron, H.H.: *Practical Subtree Transformational Grammars*; University of California at Santa Cruz, MS Thesis, 1974
- [Kr75] Kron, H.H.: *Tree Templates and Subtree Transformational Grammars*; PhD Thesis, University of California at Santa Cruz, 1975
- [Lo71] London, R. L.: *Correctness of two compilers for a LISP subset*; AIM - 151, Stanford University, 1971
- [Lo72] London, R. L.: *Correctness of a compiler for a LISP subset*; Proceedings of an ACM conference on proving assertions about programs, Sigplan Notices 7/1 1972
- [Lo78] London, R.L. et al: *Proof rules for the programming language Euclid*; Acta Informatica, Vol. 10, No 1 (1978)
- [LW69] Lucas, P., Walk, K.: *On the formal description of PL/I*; Annual Review in Automatic Programming 6, 3 (1969)
- [LS75] Luckham, D. C., Suruki, N.: *Automatic Program Verification IV: Proof of Termination Within a Weak Logic of Programs*; AIM-269, October 1975
- [LS76] Luckham, D. C., Suruki, N.: *Automatic Program Verification V: Verification-Oriented Proof Rules for Arrays, Records and Pointers*; AIM-278, March 1976
- [Ly78] Lynn, D. S.: *Interactive Compiler Proving using Hoare proof rules*; ISI/RR-78-70, 1978
- [MP73] Manna, Z., Pnueli, A.: *Axiomatic Approach to Total Correctness of Programs*; AIM-210, July 1973
- [Ma74] Manna Z.: *Mathematical Theory of Computation*; McGraw-Hill, New York 1974
- [Ma71] Mazurkiewicz, A.: *Proving Algorithms by Tail Functions*; Information and Control, 18 (1971), pp 220-226
- [Mc62] McCarthy J.: *Towards a mathematical theory of computation*; Proceedings ICIP 1962
- [Mc63] McCarthy, J.: *A basis for a mathematical theory of computation*; in P Bratford and D. Hirschberg (eds.), Computer Programming and Formal systems, pp. 33-70, North-Holland, Amsterdam, 1963
- [Mc66] McCarthy, J.: *A formal description of a subset of Algo*; Formal language description languages for computer programming (Steel, T.B., ed.), North Holland (1966)
- [Mc66] McCarthy, J., Painter, J.: *Correctness of a Compiler for Arithmetic Expressions*, AIM-40, Stanford University 1966

- [Mc65] McKeeman, W. M.: *Peephole Optimization*; CACM 8/7 (1965)
- [Me64] Mendelson, E.: *Introduction to mathematical logic*; van Nostrand, 1964
- [MG79] Meyer, A.R., Greif, I.: *Can partial correctness assertions specify programming language semantics?*; LNCS 67 (1979)
- [MS76] Milne, R., Strachey, C.: *A theory of programming language semantics*; Chapman and Hall, London 1976
- [Mi72a] Milner, R.: *Implementation and applications of Scott's Logic for Computable Functions*; in: Proceedings of the ACM Conference on Proving Assertions about Programs, Las Cruces, 1972.
- [Mi72b] Milner, R.: *Logic for Computable Functions: Description of a Machine Implementation*; AIM-169, May 1972
- [Mi72c] Milne, R.E.: *The mathematical semantics of Algo 68 (manuscript)*; Programming Research Group, University of Oxford (1972)
- [MW72] Milner, R., Weybrauch, R.: *Proving compiler correctness in a mechanized logic*; Machine Intelligence 7, 1972
- [Mo72] Morris, F. L.: *Correctness of Translation of Programming Languages, an algebraic approach*; AIM 174, 1972
- [Mo73] Morris, F.L.: *Advice on structuring compilers and proving them correct*; Proc. ACM symp. on Principles of Programming Languages, Boston, 1973
- [Mo75a] Moses, P.D.: *The Mathematical Semantics of Algo 60*; Tech. Monograph PRG-12, Programming Research Group, University of Oxford (1975)
- [Mo75b] Moses, P.D.: *Mathematical Semantics and Compiler Generation*; PhD thesis, University of Oxford (1975)
- [No78] Nelson, G., Oppen, D. C.: *Simplification by Cooperating Decision Procedures*; AIM-311, April 1978
- [Ne73] Newey, M.: *Axioms and Theorems for Integers, Lists and Finite Sets in LCF*; AIM-184, January 1973
- [Or73] Organic, E.I.: *Computer system Organization, The B5700/B6700 Series*; Academic Press, 1973
- [Pa79] Pagan, F. A.: *Algo 68 as a metalanguage for denotational semantics*; The Computer Journal, Vol 22/1 1979
- [Pa67] Painter, J. A.: *Semantic correctness of a compiler for an AlgoL-like language*; AIM - 44, Stanford University, 1967
- [Pa69] Park, D.: *Fixed point induction and proofs of programs*; Machine Intelligence 4, Edinburgh University Press, 1969
- [Pl78] Plotkin, G.: *T<sup>w</sup> as a Universal Domain*; Journal of computer and system sciences 17, pp 209-236 (1978)
- [Pl76] Plotkin G.: *A powerdomain construction*; SIAM Journal of Computing 5, 1976, 452-487.
- [Po79] Polak, W.: *An exercise in automatic program verification*; IEEE Transactions on Software Engineering, SE-5/5 (1979)
- [Po80] Polak, W.: *Program verification based on denotational semantics*; forthcoming
- [Ran64] Randell, Russell: *Algo 60 Implementation*, Academic Press, London 1964
- [Re72] Reynolds, J.C.: *Definitional Interpreters for High-order Programming languages*; Proc. 25th ACM National Conf, Boston (1972), pp 717-740
- [Re74] Reynolds, J.C.: *On the Relation between Direct and Continuation Semantics*; Proc. 2nd Coll. on Automata, languages and Programming, Saarbrücken, pp. 157 - 168, 1974
- [Re79] Reynolds J.C.: *Reasoning about arrays*; Communications of the ACM 22/5, 1979, 290-299.
- [Ro71] Rosen, B.K.: *Suttree Replacement Systems*; TR 2-71, Harvard University, 1971
- [Sa75] Samet, H.: *Automatically Proving the Correctness of Translations Involving Optimized Code*; AIM-259, May 1975
- [Sc76a] Schmeck, H.: *Ein algebraischer Ansatz für Komplierkorrektheitsbeweise*; (Brauer, W., ed.) Informatik - Fachberichte, Programmiersprachen, Springer, Berlin, Heidelberg, New York, 1976

AD-A094 604

STANFORD UNIV. CA DEPT OF COMPUTER SCIENCE  
THEORY OF COMPILER SPECIFICATION AND VERIFICATION. (U)

MAY 80 W H POLAK

F/6 9/2

MDA903-76-C-0206  
NL

UNCLASSIFIED

2 OF 2  
80246 604

END  
2-84  
DTIC

## References

- 165      166
- 
- [Sc78] Schwartz, R.L.: *An Axiomatic Semantic Definition of ALGOL 68*; CS Dept., UCLA, UCLA-34-P214-75, Aug. 78
- [Sc70] Scott, D.: *Outline of a Mathematical Theory of Computation*; Proc. 4th annual Princeton Conf. on Information Sciences and Systems, Princeton University (1970), pp169-176
- [SS71] Scott, D., Strachey, C.: *Toward a Mathematical Semantics for Computer Languages*; Tech. Monograph PRG-6, Programming Research Group, University of Oxford (1971)
- [Sc72a] Scott D.: *Continuous Lattices*; Springer Lecture Notes in Mathematics, vol. 274, 97-136, 1972.
- [Sc72b] Scott, D.: *Lattice Theory, Data Types and Semantics*; NYU Symp. on Formal Semantics, Prentice-Hall, New York (1972)
- [Sc76b] Scott, D.: *Data Types as Lattices*; SIAM Journal of Computing, 5 (1976), pp522-587
- [Sm78] Smyth M.B.: *Power domains*; Journal of Computer and System Sciences 16, 1978, 23-36.
- [SV79] Stanford Verification Group: *Stanford Pascal Verifier User Manual*; Stanford Verification Group Report No. 11, 1979
- [St77] Stoy, J.: *Denotational Semantics — The Scott-Strachey Approach to Language Theory*, MIT Press, Cambridge (1977)
- [SW74] Strachey, C., Wadsworth, C. P.: *Continuations, a Mathematical Semantics for Handling Full Jumps*; Technical Monograph PRG-11, Oxford University, 1974
- [Su75] Suzuki, N.: *Verifying Programs by Algebraic and Logical Reduction*; Proceedings of Int'l Conf on Reliable Software, IEEE, pp 473-481 (1975)
- [Su76] Suzuki, N.: *Automatic Verification of Programs with Complex Data Structures*; AIM-279, February 1976
- [Su80] Suzuki, N.: *Analysis of pointer rotation*; Seventh Annual ACM Symp. on Principles of Programming languages, Las Vegas (1980)
- [Te76] Tennent, R.D.: *The Denotational Semantics of Programming Languages*; CACM, 19 (1976) pp437-453
- [Te77a] Tennent, R.D.: *A Denotational Definition of the Programming language Pascal*; Tech. Report 77-47, Queen's University, Kingston, Ontario (1977)
- [Te77b] Tennent, R.D.: *language design methods based on semantic principles*; Acta Informatica 8, (1977)
- [TW79] Thatcher, J. W., Wagner, E. G., Wright, J.B.: *More advice on structuring compilers and proving them correct*; Report RC 7588, IBM Yorktown Heights, 1979
- [Wi69] van Wijngaarden, A. et al.: *Report on the Algorithmic language Algol 68*; Numerische Mathematik, 14, pp 79-218. 1969
- [Wi76] van Wijngaarden, A. et al.: *Revised Report on the Algorithmic language Algol 68*; Springer Berlin Heidelberg New York, 1976
- [We72] Wegner, P.: *The Vienna definition language*; Comp. Surveys, 4/1 (1972)
- [WS77] Welsh, J., Sneedinger, W.J., Hoare, C.A.R.: *Ambiguities and Insecurities in Pascal*; Software practice and experience 7 (1977)
- [WM72] Weyhrauch, R., Milner R.: *Program semantics and correctness in a mechanized logic*; Proceedings of the First USA-Japan Computer Conference, 1972.
- [WT74] Weyhrauch, R., W. Thomas, A. J.: *FOL a Proof Checker for First-order Logic*; AIM-235, 1974

**Appendix 1. Formal definition of LS****1. Micro Syntax of LS****1.1. Domains** $c \in Ch$  $s \in Str = Ch^*$  $T^m$  $V^l = N + \{\text{void}\}$  $T_k = T^m \times V^l$  $h \in H = (L_{id} \rightarrow N) \times (N \rightarrow \{\text{used}, \text{unused}\})$ 

**Identifier Tables**  
 We use  $N$  to encode token values. Operators are assigned a unique number. Identifiers are encoded as numbers. Numerals are represented by the number they denote.

**1.2. Languages  $L_i$** 

$$\begin{aligned} le &= \{"a", "b", "c", \dots, "x", "y", "z"\} \\ di &= \{"1", "2", \dots, "9", "0"\} \\ dl &= \{"+", "-", "*", "/", "<", ">", "=", ",", ";", ",", "\{", "\}", "\!", "\&"\} \\ pe &= \{"(", ")"\} \\ co &= \{".", "\;"\} \end{aligned}$$
 $L_{id} = le (le + di)^*$  $L_n = di \cdot di^*$  $L_d = de$  $L_p = pe + (pe \cdot pe)$  $L_c = co + (co \cdot co)$ **1.3. Auxiliary definitions** $\kappa \in K = L_{id} \rightarrow (T_k + \{?\})$ 

$$\begin{aligned} \kappa [[\text{array}]] &= (\text{arraysymbol}, \text{void}) \\ \kappa [[\text{begin}]] &= (\text{beginsymbol}, \text{void}) \\ \kappa [[\text{const}]] &= (\text{constsymbol}, \text{void}) \\ \kappa [[\text{do}]] &= (\text{dosymbol}, \text{void}) \\ \kappa [[\text{else}]] &\geq (\text{elsesymbol}, \text{void}) \\ \kappa [[\text{end}]] &= (\text{endsymbol}, \text{void}) \\ \kappa [[\text{fi}]] &= (\text{fisymbol}, \text{void}) \\ \kappa [[\text{function}]] &= (\text{functionsymbol}, \text{void}) \\ \kappa [[\text{goto}]] &= (\text{gosymbol}, \text{void}) \\ \kappa [[\text{if}]] &= (\text{isymbol}, \text{void}) \\ \kappa [[\text{od}]] &= (\text{odsymbol}, \text{void}) \\ \kappa [[\text{of}]] &= (\text{osymbol}, \text{void}) \\ \kappa [[\text{procedure}]] &= (\text{proceduresymbol}, \text{void}) \\ \kappa [[\text{program}]] &= (\text{programsymbol}, \text{void}) \\ \kappa [[\text{record}]] &= (\text{recordsymbol}, \text{void}) \\ \kappa [[\text{repeat}]] &= (\text{repeatsymbol}, \text{void}) \\ \kappa [[\text{then}]] &= (\text{thensymbol}, \text{void}) \\ \kappa [[\text{type}]] &= (\text{typesymbol}, \text{void}) \\ \kappa [[\text{until}]] &= (\text{untilsymbol}, \text{void}) \\ \kappa [[\text{var}]] &= (\text{varsymbol}, \text{void}) \\ \kappa [[\text{while}]] &= (\text{whitesymbol}, \text{void}) \\ \kappa z &=? \text{ for all other cases} \end{aligned}$$
 $newn \in H \rightarrow (H \times N)$ 

$$\begin{aligned} newn \cdot h &= (h[[\text{used}/n]], n) \\ \text{where } h \cdot n &= \text{unused} \end{aligned}$$

$$N \in L_n \rightarrow N$$

$$\begin{aligned} N() &= 0 \\ N.d(d_0) &= \text{if } d_0 = "0" \text{ then } 0 + 10^*(N.d) \text{ else} \\ &\quad \dots \\ &\quad \text{if } d_0 = "8" \text{ then } 8 + 10^*(N.d) \text{ else } 9 + 10^*(N.d) \end{aligned}$$

#### 1.4. Semantic Functions

$$S_{id} \in H \rightarrow L_{id} \rightarrow (Tk \times H)$$

$$\begin{aligned} S_{id}.h.s &= \text{if } s \neq ? \text{ then } (\kappa.s, h) \text{ else} \\ &\quad \text{if } h.s \neq ? \text{ then } (h.s, h) \text{ else } \langle\langle idsymbol, n\rangle, h\rangle \\ &\quad \text{where } \langle\hat{n}, n\rangle = \text{newn } h \\ &\quad \text{where } \hat{n} = \hat{h}[n/\hat{n}] \end{aligned}$$

$$S_n \in H \rightarrow L_n \rightarrow (Tk \times H)$$

$$S_n.h.s = \langle\langle numbersymbol, N.s\rangle, h\rangle$$

$$S_d \in H \rightarrow L_d \rightarrow (Tk \times H)$$

$$\begin{aligned} S_d.h.s &= \text{if } s = "+" \text{ then } \langle\langle op4symbol, 1\rangle, h\rangle \text{ else} \\ &\quad \text{if } s = "-" \text{ then } \langle\langle op4symbol, 2\rangle, h\rangle \text{ else} \\ &\quad \text{if } s = "*" \text{ then } \langle\langle op5symbol, 3\rangle, h\rangle \text{ else} \\ &\quad \dots \end{aligned}$$

Operator symbols are classified as op1symbol, op2symbol, etc. according to their precedence. The extract operator is encoded as the token value.

$$S_p \in H \rightarrow L_p \rightarrow (Tk \times H)$$

$$\begin{aligned} S_p.h.s &= \text{if } s = ":" \text{ then } \langle\langle periodsymbol, void\rangle, h\rangle \text{ else} \\ &\quad \langle\langle periodsymbol, void\rangle, h\rangle \end{aligned}$$

$$S_c \in H \rightarrow L_c \rightarrow (Tk \times H)$$

$$\begin{aligned} S_c.h.s &= \text{if } s = ":" \text{ then } \langle\langle colonsymbol, void\rangle, h\rangle \text{ else} \\ &\quad \langle\langle decompresssymbol, void\rangle, h\rangle \end{aligned}$$

$$\text{scan} \in Str \rightarrow Tk^*$$

$$\text{scan} = \Psi.h_0.$$

$$\begin{aligned} h_0.i &= \text{if } i = "\text{read}" \text{ then } 1 \text{ else} \\ &\quad \text{if } i = "\text{write}" \text{ then } 2 \text{ else} \\ &\quad \text{if } i = "\text{eof}" \text{ then } 3 \text{ else} \\ &\quad \text{if } i = "\text{integer}" \text{ then } 4 \text{ else} \\ &\quad \text{if } i = "\text{boolean}" \text{ then } 5 \text{ else} \\ &\quad \text{if } i = "\text{true}" \text{ then } 6 \text{ else} \\ &\quad \text{if } i = "\text{false}" \text{ then } 7 \text{ else } \perp \\ h_0^{\#2}.n &= \text{if } 1 \leq n \leq 7 \text{ then used else unused} \end{aligned}$$

$$\Psi \in H \rightarrow Str \rightarrow Tk^*$$

$$\begin{aligned} \Psi.h.s &= \text{if } s = () \text{ then } () \text{ else} \\ &\quad \text{if } hd.s \in \bigcup a_i \text{ then } \Phi.h(hd.s)(tl.s) \\ &\quad \Psi.h([l.s]) \\ \text{where } a_i &= \{u \mid \exists v. uv \in L_i, u \neq ()\} \end{aligned}$$

$$\Phi \in H \rightarrow Str \rightarrow Tk^*$$

$$\begin{aligned} \Phi.h.s_1.s_2 &= \text{if } s_2 = () \text{ then } (S_i.h.s_1)^{\#1} \text{ else error else} \\ &\quad \text{if } s_1.(hd.s_2) \in \bigcup a_i \text{ then } \Phi.h(s_1.(hd.s_2))(tl.s_2) \text{ else} \\ &\quad \text{if } s_1 \in L_i \text{ then } (S_i.h.s_1)^{\#1}.\Psi(S_i.h.s_1)^{\#2} \text{ else error} \end{aligned}$$

## 2. Syntax of LS

The following SLR-grammar defines the syntax of LS. All uppercase symbols are terminals, lower case symbols are nonterminals.

[Z..1]	Z ::= PROG eofsymbol	PROC ::= programsymbol STMT	COM ::= COM semicolon symbol LST
[PROG..1]		BLOCK ::= CDEFP TDEF CP PFDEC CP CSTMT	STMT ::= STMT
[CDEFP..1]	CDEFP ::= CDEFL	CDEFP ::= CDEFL CDEF P VDEF CP VFDEC CP	STMT ::= VBLE becomesymbol EXPR
[CDEFP..2]	CDEF ::= constsymbol	CDEF ::= constsymbol CDEF L	STMT ::= idsymbol
[CDEFL..1]	CDEFL ::= CDEF	CDEFL ::= CDEF semicolon symbol	STMT ::= idsymbol lparentsymbol EXPR lparentsymbol
[CDEFL..2]	CDEF ::= equalsymbol	CDEF ::= CDEF semicolon symbol EXPR	STMT ::= idsymbol ifsymbol EXPR thensymbol COM elsesymbol COM fiymbol
[CDEF..1]	CDEF ::= idsymbol	CDEF ::= idsymbol equalsymbol EXPR	STMT ::= idsymbol EXPR thensymbol COM dosymbol COM odymbol
[TDEF..1]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF L	STMT ::= whitespace symbol EXPR
[TDEF..2]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	STMT ::= repeatsymbol COM utilitysymbol EXPR
[TDEF..3]	TDEF ::= TDEF	TDEF ::= TDEF TDEF semicolon symbol	STMT ::= gotosymbol COM numbersymbol
[TDEF..4]	TDEF ::= idsymbol	TDEF ::= idsymbol TDEF	STMT ::= BLOCK
[TDEF..5]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	STMT ::= EXPR op1symbol CONJ
[TDEF..6]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	STMT ::= EXPR op2symbol CONJ
[TDEF..7]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	STMT ::= CONJ CONJ
[TDEF..8]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	CONJ ::= CONJ op2symbol REL
[TDEF..9]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	CONJ ::= CONJ REL
[TDEF..10]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	CONJ ::= unopsymbol REL
[TDEF..11]	TDEF ::= typesymbol	TDEF ::= typesymbol TDEF semicolon symbol	REL ::= REL op3symbol SUM
[TDEN..1]	TDEN ::= lparentsymbol IDL rparentsymbol	TDEN ::= lparentsymbol IDL rparentsymbol	REL ::= REL equalsymbol SUM
[TDEN..2]	TDEN ::= idsymbol	TDEN ::= idsymbol	REL ::= SUM
[TDEN..3]	TDEN ::= idsymbol colonsymbol IDN	TDEN ::= idsymbol colonsymbol IDN	SUM ::= SUM op4symbol TERM
[TDEN..4]	TDEN ::= numbersymbol colonsymbol IDN	TDEN ::= numbersymbol colonsymbol IDN	SUM ::= TERM
[TDEN..5]	TDEN ::= upsymbol idsymbol	TDEN ::= upsymbol idsymbol	SUM ::= op5symbol TERM
[TDEN..6]	TDEN ::= arraysymbol bracketsymbol TDEN bracketsymbol of symbol TDEN	TDEN ::= arraysymbol bracketsymbol TDEN bracketsymbol of symbol TDEN	TERM ::= TERM op5symbol FACT
[TDEN..7]	TDEN ::= recordsymbol VDECL endsymbol	TDEN ::= recordsymbol VDECL endsymbol	TERM ::= FACT TERM
[IDL..1]	IDL ::= idsymbol	IDL ::= idsymbol	FACT ::= FACT
[IDL..2]	IDL ::= IDL commasymbol idsymbol	IDL ::= IDL commasymbol idsymbol	FACT ::= idsymbol lparentsymbol EXPR rparentsymbol
[VDEC..1]	VDEC ::= VDEC	VDEC ::= VDEC	FACT ::= IDN
[VDEC..2]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	FACT ::= FACT upsymbol
[VDEC..3]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	FACT ::= FACT periodsymbol idsymbol
[VDEC..4]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	FACT ::= FACT lbracketsymbol EXPR rbracketsymbol
[VDEC..5]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	VBLE ::= IDN
[VDEC..6]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	VBLE ::= VBLE
[VDEC..7]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	VBLE ::= VBLE upsymbol
[VDEC..8]	VDEC ::= VDEC	VDEC ::= VDEC semicolon symbol VDEC	VBLE ::= VBLE periodsymbol idsymbol
[PFDEC..1]	PFDEC ::= PFDEC	PFDEC ::= PFDEC	VBLE ::= VBLE lbracketsymbol EXPR rbracketsymbol
[PFDEC..2]	PFDEC ::= PFDEC	PFDEC ::= PFDEC semicolon symbol PFDEC	PARMS ::= PARM
[PFDEC..3]	PFDEC ::= PFDEC	PFDEC ::= PFDEC semicolon symbol PFDEC	PARMS ::= iparensymbol PARM
[FDEC..1]	FDEC ::= functionsymbol idsymbol PARMS colonsymbol TDEN semicolon symbol STMT	FDEC ::= functionsymbol idsymbol PARMS colonsymbol TDEN semicolon symbol STMT	PARML ::= PARM
[CSTM..1]	CSTM ::= proceduresymbol idsymbol PARMS semicolon symbol STMT	CSTM ::= proceduresymbol idsymbol PARMS semicolon symbol STMT	PARML ::= PARM
[LST..1]	LST ::= beginsymbol COM endsymbol	LST ::= beginsymbol COM endsymbol	PARM ::= varsymbol idsymbol colonsymbol idsymbol
[LST..2]	LST ::= numbersymbol colonsymbol STMT	LST ::= numbersymbol colonsymbol STMT	PARM ::= idsymbol colonsymbol idsymbol
[COM..1]	COM ::= LST	COM ::= LST	EXPR ::= EXPR
			EXPR ::= EXPR commasymbol EXPR

Appendix 1. Formal definition of LS

$|IDN_1|$        $IDN ::= \text{idsymbol}$   
 $|IDN_2|$        $IDN ::= \text{number symbol}$

### 3. Abstract syntax

#### 3.1. Syntactic Domains

$\Omega \in Op$                           dyadic Operators (not further defined here)  
 $O \in Mop$                           monadic Operators (not further defined here)  
 $I \in Id$                                   Identifiers (not further defined here)  
 $N \in Num$                                   Numerals (not further defined here)  
 $B \in Pgm$                                   Programs  
 $E \in Exp$                                   Expressions  
 $\Theta \in SIm$                                   Statements  
 $\Gamma \in Com$                                   Commands  
 $T \in Typ$                                   Types  
 $\Delta \in Dec$                                   Declarations  
 $\Delta_c \in Cdef$                                   Constant Definition  
 $\Delta_t \in Tdef$                                   Type definition  
 $\Delta_v \in Vdec$                                   Variable declaration  
 $\Pi \in Par$                                   Parameters

**Expressions:**  
 $E ::= I \mid O E \mid E_0 \Omega E_1 \mid E \uparrow \mid I(E^*) \mid E_0[E_i] \mid N \mid E \cdot I$

**Commands:**  
 $\Gamma ::= N : \Theta \mid \Theta \mid \Gamma_0 ; \Gamma_1$

**Statements:**

$\Theta ::= E_0 ::= E_1 \mid \text{if } E \text{ then } \Gamma_0 \text{ else } \Gamma_1 \text{ fi } \mid \text{dummy} \mid$   
 $\text{while } E \text{ do } \Gamma \text{ od } \mid \text{repeat } \Gamma \text{ until } E \mid \text{goto } N \mid I(E^*) \mid$   
 $\Delta^* \begin{array}{l} \text{begin } \Gamma \text{ end} \end{array}$

**Declarations:**

$\Delta ::= \text{const } \Delta^* \mid \text{type } \Delta^* \mid \text{var } \Delta^* \mid$   
 $\text{procedure } I(\Pi^*) ; \Theta \mid \text{function } I(\Pi^*) ; \Theta$

**Types:**

$T ::= I \mid (I_1, \dots, I_n) \mid E_1 \cdot E_2 \mid \text{array}[T_1]_0 / T_2 \mid$

$\text{record } I_1; T_1; I_2; T_2; \dots; I_n; T_n \text{ end } \mid I \uparrow I$

**Constant Definitions:**  
 $\Delta_c ::= I = E$

**Type Definitions:**

$\Delta_t ::= I = T$

**Variable Declarations:**

$\Delta_v ::= I \cdot T$

**Parameters:**  
 $\Pi ::= I_1; I_2 \mid \text{var } I_1; I_2$

#### 3.2. Constructor functions

**Statements:**  
 $mkeif(E, \Gamma_1, \Gamma_2) = \text{if } E \text{ then } \Gamma_1 \text{ else } \Gamma_2 \text{ fi}$   
 $mkeassign(E_1, E_2) = E_1 := E_2$   
 $mkewhile(E, \Gamma) = \text{while } E \text{ do } \Gamma \text{ od}$   
 $mkerpeat(\Gamma, E) = \text{repeat } \Gamma \text{ until } E$   
 $mkegoto(N) = \text{goto } N$   
 $mkecall(I, (E_1, \dots, E_n)) = I(E_1, \dots, E_n)$   
 $mkeblock((\Delta_1, \dots, \Delta_n), \Gamma) = \Delta_1, \dots, \Delta_n \begin{array}{l} \text{begin } \Gamma \text{ end} \end{array}$   
 $mkedummy = \text{empty}$

**Expressions:**

$mkeexpid(I) = I$   
 $mkenumber(N) = E$   
 $mkeunop(O, E) = OE$   
 $mkebinop(E_1, 1, E_2) = E_1 + E_2$   
 $mkebinop(E_1, 2, E_2) = E_1 - E_2$   
 $\dots$   
 $mkederef(E) = E \uparrow$   
 $mkefcall(I, (E_1, \dots, E_n)) = I(E_1, \dots, E_n)$   
 $mkeindex(E_1, E_2) = E_1[E_2]$   
 $mkeselect(E, I) = E \cdot I$

**Commands:**

$mklabel(N, \Theta) = N : \Theta$   
 $mlestmi(\Theta) = \Theta$

$mkecommandlist(\Gamma_1, \Gamma_2) = \Gamma_1; \Gamma_2$ 

Declarations:

$mktyp((\Delta_{c1}, \dots, \Delta_{cn})) = type \Delta_{c1}, \dots, \Delta_{cn}$   
 $mkeconst((\Delta_{c1}, \dots, \Delta_{cn})) = const \Delta_{c1}, \dots, \Delta_{cn}$

$mkevar((\Delta_{v1}, \dots, \Delta_{vn})) = var \Delta_{v1}, \dots, \Delta_{vn}$   
 $mkefunction(I, (\Pi_1, \dots, \Pi_n), T, \Theta) = function I(\Pi_1, \dots, \Pi_n); T; \Theta$   
 $mkeprocedure(I, (\Pi_1, \dots, \Pi_n), \Theta) = procedure I(\Pi_1, \dots, \Pi_n); \Theta$

Parameters:

$mkevarp(I_1; I_2) = var I_1; I_2$   
 $mkevalp(I_1; I_2) = I_1; I_2$

Types:

$mkeenum((I_1, \dots, I_n)) = (I_1, \dots, I_n)$   
 $mktypeset(I) = I$   
 $mkearray(T_1, T_2) = array [T_1] of T_2$   
 $mkerecord(I; T_1, \dots, I_n; T_n) = record I; T_1, \dots, I_n; T_n end$   
 $mkesubrange(E_1, E_2) = E_1 \dots E_1$   
 $mkepointer(T) = \uparrow T$

Programs

 $mkeprogram(\Theta) = program \Theta$ .

In addition to these constructors there is a set of list constructors which are used to represent lists like  $E_1, \dots, E_n$ .

$mkeelist$  — make a singleton list  
 $mkeappend$  — append two lists  
 $mkenullist$  — make an empty list  
 $first$  — true for lists  
 $isnullist$  — true for empty lists  
 $selfirst$  — select the first element  
 $selrest$  — select the rest

 $mkearray$  — make an array $mkepointer$  — make a pointer $mkesubrange$  — make a subrange $mkepointer$  — make a record $mkearray$  — make an array $mkepointer$  — make a pointer $mkearray$  — make an array $mkepointer$  — make a pointer $mkearray$  — make an array $mkepointer$  — make a pointer $mkearray$  — make an array $mkepointer$  — make a pointer**4. Tree transformations****4.1. Programs** $\mathcal{E}(t, v) = v$ 

$\mathcal{E}(z, I, r_1, r_2) = \mathcal{E}(r_1)$   
 $\mathcal{E}(\text{PROG}, I, r_1, r_2) = mkeprogram(\mathcal{E}(r_2))$   
 $\mathcal{E}(\text{BLOCK}, I, r_1, r_2, r_3, r_4, r_5) =$   
 $mkeblock(mkeappend(\mathcal{E}(r_1), \mathcal{E}(r_2), \mathcal{E}(r_3), \mathcal{E}(r_4)), \mathcal{E}(r_5))$

**4.2. Declarations**

$\mathcal{E}(\text{CDEFP}, 1) = mkenullist$   
 $\mathcal{E}(\text{CDEFP}, 2, r_1, r_2) = mkeconst(\mathcal{E}(r_2))$   
 $\mathcal{E}(\text{CDEFL}, I, r_1, r_2) = mkeconst(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{CDEFL}, 2, r_1, r_2, r_3) = mkeappend(\mathcal{E}(r_1), mkeconst(\mathcal{E}(r_2)))$   
 $\mathcal{E}(\text{CDEF}, I, r_1, r_2, r_3) = mkeconstdec(\mathcal{E}(r_1), \mathcal{E}(r_3))$   
 $\mathcal{E}(\text{TDEFP}, 1) = mkenullist$   
 $\mathcal{E}(\text{TDEFP}, 2, r_1, r_2, r_3) = mkectype(\mathcal{E}(r_2))$   
 $\mathcal{E}(\text{TDDEFL}, 1, r_1, r_2) = mkeconst(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{TDDEFL}, 2, r_1, r_2, r_3) = mkeappend(\mathcal{E}(r_1), mkeconst(\mathcal{E}(r_2)))$   
 $\mathcal{E}(\text{TDEF}, I, r_1, r_2, r_3) = mkeypecad(\mathcal{E}(r_1), \mathcal{E}(r_3))$   
 $\mathcal{E}(\text{TDEF}, 2, r_1, r_2, r_3) = mkeconst(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{TDEN}, I, r_1, r_2, r_3) = mkeenum(\mathcal{E}(r_2))$   
 $\mathcal{E}(\text{TDEN}, 2, r_1) = mketypdec(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{TDEN}, 3, r_1, r_2, r_3) = mkesubrange(\mathcal{E}(r_1), \mathcal{E}(r_3))$   
 $\mathcal{E}(\text{TDEN}, 4, r_1, r_2, r_3) = mkesubrange(\mathcal{E}(r_1), \mathcal{E}(r_3))$   
 $\mathcal{E}(\text{TDEN}, 5, r_1, r_2) = mkepointer(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{TDEN}, 6, r_1, r_2, r_3, r_4, r_5, r_6) = mkearray(\mathcal{E}(r_3), \mathcal{E}(r_6))$   
 $\mathcal{E}(\text{TDEN}, 7, r_1, r_2, r_3) = mkerecord(\mathcal{E}(r_3))$   
 $\mathcal{E}(\text{IDL}, I, r_1) = mkeconst(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{IDL}, 2, r_1, r_2, r_3) = mkeappend(\mathcal{E}(r_1), mkeconst(\mathcal{E}(r_3)))$   
 $\mathcal{E}(\text{VDECL}, I, r_1) = mkeconst(\mathcal{E}(r_1))$   
 $\mathcal{E}(\text{VDECL}, 2, r_1, r_2, r_3) = mkeappend(\mathcal{E}(r_1), mkeconst(\mathcal{E}(r_3)))$   
 $\mathcal{E}(\text{VDEC}, I, r_1, r_2, r_3) = mkevardcl(\mathcal{E}(r_1), \mathcal{E}(r_3))$   
 $\mathcal{E}(\text{VDEC}, 1) = mkenullist$   
 $\mathcal{E}(\text{VDEC}, 2, r_1, r_2, r_3) = mkevar(\mathcal{E}(r_2))$

The obvious axioms for lists hold for these functions.

We use *selfirst* and *selrest* here to indicate that these are special destructors of the abstract syntax as opposed to *hd* and *tl* which operate on arbitrary lists.

$E(PDEC.P.1) = mkenullist$	$E(EXPR.1, r_1, r_2, r_3) = mkebinop(E(r_1), E(r_2), E(r_3))$
$E(PFDEC.P.2, r_1, r_2, r_3) = mkeappend(mkeelist(E(r_1)), E(r_3))$	$E(EXPR.2, r_1) = E(r_1)$
$E(PFDEC.P.3, r_1, r_2, r_3) = mkeappend(mkeelist(E(r_1)), E(r_3))$	$E(CONJ.1, r_1, r_2, r_3) = mkebinop(E(r_1), E(r_2), E(r_3))$
$E(FDEC.1, r_1, r_2, r_3, r_4, r_5, r_6, r_7) = mkefunction([E(r_2)], E(r_3), E(r_5), E(r_7))$	$E(CONJ.2, r_1) = E(r_1)$
$E(PDEC.1, r_1, r_2, r_3, r_4, r_5, r_6, r_7) = mkeprocedure([E(r_2)], E(r_3), E(r_5))$	$E(REL.1, r_1, r_2, r_3) = mkeunop(E(r_1), E(r_2), E(r_3))$
	$E(REL.2, r_1, r_2, r_3) = mkebinop(E(r_1), E(r_2), E(r_3))$
	$E(REL.3, r_1) = E(r_1)$
	$E(SUM.1, r_1, r_2, r_3) = mkeunop(E(r_1), E(r_2), E(r_3))$
	$E(SUM.2, r_1) = E(r_1)$
	$E(SUM.3, r_1, r_2) = mkeunop(E(r_1), E(r_2), E(r_3))$
	$E(TERM.1, r_1, r_2, r_3) = mkebinop(E(r_1), E(r_2), E(r_3))$
	$E(TERM.2, r_1) = E(r_1)$
	$E(FACT.1, r_1, r_2, r_3) = E(r_2)$
	$E(FACT.2, r_1, r_2, r_3, r_4) = mkecall(E(r_1), E(r_3))$
	$E(FACT.3, r_1) = E(r_1)$
	$E(FACT.4, r_1, r_2) = mkederef(E(r_1))$
	$E(FACT.5, r_1, r_2, r_3) = mkeselect(E(r_1), E(r_3))$
	$E(FACT.6, r_1, r_2, r_3, r_4) = mkeindex(E(r_1), E(r_3))$
$E(VBLE.1, r_1) = E(r_1)$	$E(STMT.1, r_1, r_2, r_3) = E(r_3)$
$E(VBLE.2, r_1, r_2) = mkederef(E(r_1))$	$E(STMT.2, r_1) = E(r_1)$
$E(VBLE.3, r_1, r_2, r_3) = mkeselect(E(r_1), E(r_3))$	$E(STMT.3, r_1) = E(r_2)$
$E(VBLE.4, r_1, r_2, r_3, r_4) = mkeindex(E(r_1), E(r_3))$	$E(STMT.4, r_1, r_2, r_3) = mkecall(E(r_1), E(r_3))$
$E(ARMS.1) = mkenullist$	$E(LST.1, r_1, r_2, r_3) = mkelatch(E(r_1), E(r_3))$
$E(PARMS.2, r_1, r_2, r_3) = E(r_2)$	$E(LST.2, r_1) = mkestmt(E(r_1))$
$E(PARML.1, r_1) = mkelist(E(r_1))$	$E(COM.1, r_1) = E(r_1)$
$E(PARML.2, r_1, r_2, r_3) = mkeappend(E(r_1), mkelist(E(r_3)))$	$E(COM.2, r_1, r_2, r_3) = mkecommandist(E(r_1), E(r_3))$
$E(PARM.1, r_1, r_2, r_3, r_4) = mkeunop(E(r_2), E(r_4))$	$E(STMT.5, r_1, r_2, r_3, r_4, r_5, r_6, r_7) = mkestmt(mkedurnmy)$
$E(PARM.2, r_1, r_2, r_3) = mkevalp(E(r_1), E(r_3))$	$E(STMT.6, r_1, r_2, r_3, r_4, r_5, r_6, r_7) = mkef(E(r_2), E(r_4), mkestmt(mkedurnmy))$
$E(EXPR.1, r_1) = mkelist(E(r_1))$	$E(STMT.7, r_1, r_2, r_3, r_4, r_5, r_6, r_7) = mkefh(E(r_2), E(r_4), E(r_6))$
$E(EXPR.2, r_1, r_2, r_3) = mkeappend(E(r_1), mkelist(E(r_3)))$	$E(STMT.8, r_1, r_2, r_3, r_4, r_5, r_6, r_7) = mkehw(E(r_2), E(r_4), E(r_6))$
$E(DN.1, r_1) = E(r_1)$	$E(STMT.9, r_1, r_2) = mkegoto(E(r_1))$
$E(DN.2, r_1) = E(r_1)$	$E(STMT.10, r_1) = E(r_1)$

## 4.4. Statements

## 4.3. Expressions



$$\text{isreturnable} \in Ty \rightarrow Ty \rightarrow T$$

$$\text{isreturnable } r = \text{isindex } r \vee r::[\nu: \uparrow r'] \vee r::[\text{nil}]$$

$$\text{overlap} \in Ty \rightarrow Ty \rightarrow T$$

$$\begin{aligned} \text{overlap } [\nu:\text{sub}:i_1:i_2][\nu:\text{sub}:i_3:i_4] &= i_1 \leq i_4 \wedge i_3 \leq i_2 \\ \text{contains } [\nu:\text{sub}:i_1:i_2][\nu:\text{sub}:i_3:i_4] &= i_1 \leq i_3 \wedge i_4 \leq i_2 \end{aligned}$$

$$\text{union} \in Ty \rightarrow Ty \rightarrow Ty$$

$$\text{union } r_1r_2 = \text{if } r_1::[\nu:\text{sub}:i_1:i_2] \text{ then } r_2::[\nu:\text{sub}:i_3:i_4] \text{ then } [\nu:\text{sub}:i_1:i_4]$$

$$\text{type} \in M\alpha \rightarrow Ty$$

$$\begin{aligned} \text{type } [\text{var}:r] &= r \\ \text{type } [\text{val}:r] &= r \\ \text{type } [\text{const}:r] &= r \\ \text{type } [\text{ofun}:\mu_1,\dots,\mu_n;r] &= r \\ \text{type } [\text{pfun}:\mu_1,\dots,\mu_n;r] &= r \end{aligned}$$

$$\text{equal} \in Ty \rightarrow Ty \rightarrow T$$

$$\begin{aligned} \text{equal } r_1r_2 &= \text{if } r_1::[\nu:\text{sub}:i_1:i_2] \text{ then } \\ &\quad \text{if } r_2::[\nu:\text{sub}:i_1:i_2] \text{ then true else false} \\ &\quad \text{else type tag } r_1 = \text{type tag } r_2 \end{aligned}$$

$$\text{compatible} \in Ty \rightarrow Ty \rightarrow T$$

$$\begin{aligned} \text{compatible } r_1r_2 &= \text{equal } r_1r_2 \vee \text{overlap } r_1r_2 \vee \\ &\quad (r_1::[\nu: \uparrow r'] \wedge r_2 = [\text{nil}]) \end{aligned}$$

$$\text{assignable} \in M\alpha \rightarrow M\alpha \rightarrow T$$

$$\begin{aligned} \text{assignable } \mu_1\mu_2 &= \text{isvar } \mu_1 \wedge \\ &\quad \text{compatible}(\text{type } \mu_1)(\text{type } \mu_2) \wedge \\ &\quad \text{returnable}(\text{type } \mu_2) \end{aligned}$$

$$\text{passable} \in M\alpha \rightarrow M\alpha \rightarrow T$$

$$\text{passable } \mu_1\mu_2 = \text{if isvar } \mu_1 \text{ then isvar } \mu_2 \wedge \text{equal}(\text{type } \mu_1)(\text{type } \mu_2) \text{ else}$$

$$\begin{aligned} \text{isvar}, \text{isval}, \in M\alpha \rightarrow T \\ \text{isvar } \mu = \mu::[\text{var}:r] \vee \mu::[\text{ap}:r] \vee \mu::[\text{ofun}:\mu_1\dots\mu_n;r] \\ \text{isval } \mu = \mu::[\text{val}:r] \vee \mu::[\text{const}:r] \end{aligned}$$

$$\text{isbool}, \text{isint} \in Ty \rightarrow T$$

$$\begin{aligned} \text{isint } r &= (\text{type tag } r = \text{int}) \\ \text{isbool } r &= (\text{type tag } r = \text{bool}) \end{aligned}$$

$$\text{sp} \in Id \rightarrow M\alpha^* \rightarrow T$$

$$\begin{aligned} \text{sp}[\text{print}](\mu) &= \text{passable } \mu \text{ [val:int; sub:int\_}\infty\text{..int\_}\infty\text{]} \\ \text{sp}[\text{new}](\mu) &= \mu::[\text{var}:[\nu: \uparrow r]] \end{aligned}$$

$$\text{sf} \in Id \rightarrow M\alpha^* \rightarrow M\alpha$$

$$\begin{aligned} \text{sf}[\text{eof}](\nu) &= [\text{val}:[\text{bool}; \text{sub}:0:1]] \\ \text{sf}[\text{read}](\nu) &= [\text{val}:[\text{int}; \text{sub:int\_}\infty\text{..int\_}\infty]] \end{aligned}$$

$\boxed{\text{true mode}, \text{false mode} \in Md}$

$\text{true mode} = [\text{const}:0:[\text{bool}:sub:0:0]]$   
 $\text{false mode} = [\text{const}:1:[\text{bool}:sub:1:1]]$

$\boxed{\text{int const} \in I \rightarrow Md}$

$\text{int const } t = [\text{const}:i:[\text{int}:sub:t:i]]$

$\boxed{\text{integer} \in Ty}$

$\text{integer} = [\text{int}:sub:\text{int} - \infty : \text{int} + \infty]$

$\boxed{\text{boolean} \in Ty}$

$\text{boolean} = [\text{bool}:sub:0:1]$

$\boxed{\text{mkeval mode} \in Md \rightarrow Md}$

$\text{mkeval mode } \mu = \text{if isval } \mu \text{ then } \mu \text{ else } [\text{val:type } \mu]$

$\boxed{o \in O \rightarrow Md \rightarrow Md}$

$o[\neg]\mu = \text{if isbool } \mu \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t:r] \text{ then}$   
 $\quad \quad \text{if } t = 0 \text{ then false mode } \mu_2 \text{ else true mode}$   
 $\quad \quad \text{else } [\text{val:boolean}]$   
 $\quad \text{else } \perp$

$\boxed{w \in \Omega \rightarrow Md \rightarrow Md \rightarrow Md}$

$w[\wedge]\mu_1\mu_2 = \text{if isbool } \mu_1 \wedge \text{isbool } \mu_2 \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t:r_1] \text{ then}$   
 $\quad \quad \text{if } t = 0 \text{ then mkeval mode } \mu_2 \text{ else false mode}$   
 $\quad \quad \text{else if } \mu_2:[\text{const}:t:r_2] \text{ then}$   
 $\quad \quad \quad \text{if } t = 0 \text{ then mkeval mode } \mu_1 \text{ else false mode}$   
 $\quad \quad \quad \text{else } [\text{val:boolean}]$   
 $\quad \text{else } \perp$

$w[\vee]\mu_1\mu_2 = \text{if isbool } \mu_1 \wedge \text{isbool } \mu_2 \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t:r_1] \text{ then}$   
 $\quad \quad \text{if } t = 0 \text{ then true mode } \text{mkeval mode } \mu_2$   
 $\quad \quad \text{else if } \mu_2:[\text{const}:t:r_2] \text{ then}$   
 $\quad \quad \quad \text{if } t = 0 \text{ then true mode } \text{mkeval mode } \mu_1$   
 $\quad \quad \quad \text{else } [\text{val:boolean}]$   
 $\quad \text{else } \perp$

$w[+]\mu_1\mu_2 = \text{if isint } \mu_1 \wedge \text{isint } \mu_2 \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t_1:r_1] \wedge \mu_2:[\text{const}:t_2:r_2] \text{ then}$   
 $\quad \quad \text{int const}(t_1 + t_2)$   
 $\quad \text{else } [\text{val:integer}]$   
 $\quad \text{else } \perp$

$w[-]\mu_1\mu_2 = \text{if isint } \mu_1 \wedge \text{isint } \mu_2 \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t_1:r_1] \wedge \mu_2:[\text{const}:t_2:r_2] \text{ then}$   
 $\quad \quad \text{int const}(t_1 - t_2)$   
 $\quad \text{else } [\text{val:integer}]$   
 $\quad \text{else } \perp$

$w[*]\mu_1\mu_2 = \text{if isint } \mu_1 \wedge \text{isint } \mu_2 \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t_1:r_1] \wedge \mu_2:[\text{const}:t_2:r_2] \text{ then}$   
 $\quad \quad \text{int const}(t_1 * t_2)$   
 $\quad \text{else } [\text{val:integer}]$   
 $\quad \text{else } \perp$

$w[<]\mu_1\mu_2 = \text{if isint } \mu_1 \wedge \text{isint } \mu_2 \text{ then}$   
 $\quad \text{if } \mu_1:[\text{const}:t_1:r_1] \wedge \mu_2:[\text{const}:t_2:r_2] \text{ then}$   
 $\quad \quad \text{bool const}(t_1 < t_2)$   
 $\quad \quad \text{else } [\text{val:boolean}]$   
 $\quad \text{else } \perp$

$w \llbracket > \mu_1 \mu_2 = \text{if } \text{isint } \mu_1 \wedge \text{isint } \mu_2 \text{ then}$   
      $\quad \text{if } \mu_1 :: [\text{const}; \tau_1] \wedge \mu_2 :: [\text{const}; \tau_2] \text{ then}$   
          $\quad \quad \text{boolean}(\epsilon_1 > \epsilon_2)$   
      $\quad \text{else } [\text{val}; \text{boolean}]$   
      $\quad \text{else } \perp$

## 6.2. Declarations

 $d \in Dec \rightarrow U_e \rightarrow U_e$ 

$d \llbracket const[\Delta_{c1}; \dots; \Delta_{cn}] \varsigma = dc^* [\Delta_{c1}; \dots; \Delta_{cn}] \varsigma$   
 $d \llbracket type[\Delta_{t1}; \dots; \Delta_{tn}] \varsigma = fit(dt^* [\Delta_{t1}; \dots; \Delta_{tn}]) \varsigma$   
 $d \llbracket var[\Delta_{v1}; \dots; \Delta_{vn}] \varsigma = dv^* [\Delta_{v1}; \dots; \Delta_{vn}] \varsigma$

$d \llbracket procedure I(\Pi_1, \dots, \Pi_n); \Theta \rrbracket \varsigma =$   
      $\quad \text{let } ((\mu_1, \dots, \mu_n), (\varsigma_1, \dots, \Pi_n)) \in \Theta \text{ in}$   
      $\quad \quad \text{if } \text{distinct}((\Pi_1, \dots, \Pi_n); \Theta) \text{ then}$   
      $\quad \quad \quad g[\Theta](\varsigma_1[\text{proc}(\mu_1, \dots, \mu_n)] / I)$   
      $\quad \quad \quad \text{to } [\text{proc}(\mu_1, \dots, \mu_n)] / I$

$d \llbracket function I(\Pi_1, \dots, \Pi_n).I; \Theta \rrbracket \varsigma =$   
      $\quad \text{let } ((\mu_1, \dots, \mu_n), (\varsigma_1, \dots, \Pi_n)) \in \Theta \text{ in}$   
      $\quad \quad \text{if } ([\text{type } \tau] :: \varsigma_1[I_2]) \text{ then}$   
      $\quad \quad \quad \text{if } \text{returnable}(\tau) \text{ then}$   
      $\quad \quad \quad \text{if } \text{distinct}((\Pi_1, \dots, \Pi_n); \Theta) \text{ then}$   
      $\quad \quad \quad \quad \text{if } g[\Theta](\varsigma_1[\text{afun}(\mu_1, \dots, \mu_n); \tau] / I) \text{ then}$   
      $\quad \quad \quad \quad \quad \text{to } [\text{afun}(\mu_1, \dots, \mu_n); \tau] / I$

 $dt \in Tdef \rightarrow U_e \rightarrow U_e$ 

$p \llbracket I_1; I_2 \rrbracket \varsigma = \text{if } [\text{type } \tau] :: \varsigma[I_2] \text{ then}$   
      $\quad \quad \quad \text{if } \text{returnable } \tau \text{ then } ([var;\tau], \varsigma[[var;\tau]/I_1])$   
 $p \llbracket var I_1; I_2 \rrbracket \varsigma = \text{if } [type \tau] :: \varsigma[I_2] \text{ then } ([var;\tau], \varsigma[[var;\tau]/I_1])$   
 $dt[I = T] \varsigma = \text{let } (\tau, \varsigma) = t[T] \varsigma \text{ in } \varsigma[[type;\tau]/I]$

 $dt^* \in Tdf^* \rightarrow U_e \rightarrow U_e$ 

$dt^* \llbracket \text{so}_1 = \varsigma_0$   
 $dt^* \llbracket \text{so}_1, \dots, \text{so}_n = dt^* [\Delta_{s1}; \dots; \Delta_{sn}] (dt[\Delta_{s0}]\varsigma_0)\varsigma_1$

 $dc \in Cdef \rightarrow U_e \rightarrow U_e$  $dc[I = E] \varsigma = \text{if } [\text{const}; \tau] :: E \varsigma \text{ then } \varsigma[[\text{const}; \tau]/I]$  $dc^* \in Cdef^* \rightarrow U_e \rightarrow U_e$ 

$dc^* \llbracket \varsigma = \varsigma$   
 $dc^* \llbracket \Delta_{c0}; \Delta_{c1}; \dots; \Delta_{cn} \varsigma = dc^* [\Delta_{c1}; \dots; \Delta_{cn}] (dc[\Delta_{c0}]\varsigma)$

 $dv \in Vdec \rightarrow U_e \rightarrow U_e$  $dv \llbracket I; T \rrbracket \varsigma = \text{let } (\tau, \varsigma) = t[T] \varsigma \text{ in } \varsigma[[var;\tau]/I]$  $dv^* \in Vdec^* \rightarrow U_e \rightarrow U_e$ 

$dv^* \llbracket \varsigma = \varsigma$   
 $dv^* \llbracket \Delta_{v0}; \Delta_{v1}; \dots; \Delta_{vn} \varsigma = dv^* [\Delta_{v1}; \dots; \Delta_{vn}] (dv[\Delta_{v0}]\varsigma)$

 $p \in Par \rightarrow U_e \rightarrow (Md \times U_e)$ 

$p \llbracket I_1; I_2 \rrbracket \varsigma = \text{if } [\text{type } \tau] :: \varsigma[I_2] \text{ then}$   
      $\quad \quad \quad \text{if } \text{returnable } \tau \text{ then } ([var;\tau], \varsigma[[var;\tau]/I_1])$   
 $p \llbracket var I_1; I_2 \rrbracket \varsigma = \text{if } [type \tau] :: \varsigma[I_2] \text{ then } ([var;\tau], \varsigma[[var;\tau]/I_1])$

$$P^* \in Par^* \rightarrow U_* \rightarrow (Md^* \times U_*)$$

$$P^*[\pi_1, \dots, \pi_n]_{\varsigma_0} = ((\mu_1, \dots, \mu_n), \varsigma_n)$$

where  $(\mu_i, \varsigma_i) = P[\pi_i]_{\varsigma_{i-1}}$

$$d^* \in Dec^* \rightarrow U_* \rightarrow U_*$$

$$d^*[(\Delta_1, \dots, \Delta_n)]_{\varsigma} = d^*[(\Delta_2, \dots, \Delta_n)](d[\Delta_1]_{\varsigma})$$

$$t \in Typ \rightarrow U_* \rightarrow U_* \rightarrow (Ty \times U_*)$$

$$t[I]_{\varsigma_0} = \text{if } [type\ r]::[I] \text{ then } (r, \varsigma_0)$$

$$t[(I_1, \dots, I_n)]_{\varsigma_0} = ((\nu, sub:I:n), \varsigma_2[\mu_i/I_i])$$

$$\text{where } \mu_i = [\text{const}:[\nu:i:i]]$$

$$\text{where } (\nu, \varsigma_2) = \text{newtag } \varsigma_0$$

$$t[E_1 \dots E_2]_{\varsigma_0} = \text{if } [\text{const}:(\tau_1 \dots e[E_1])_{\varsigma_0} \text{ then}$$

$$((\text{union } \tau_1 \tau_2), \varsigma_0)$$

$$t[\text{array}[T_1]_1, T_2]_{\varsigma_0} = \text{let } (r_1, \varsigma_1) = t[T_1]_{\varsigma_0}, (r_2, \varsigma_2) = t[T_2]_{\varsigma_2} \text{ in}$$

if isindex  $r_1$  then  $((\nu:\text{array}:\tau_1:\tau_2), \varsigma_1)$

$$\text{where } (\nu, \varsigma_2) = \text{newtag } \varsigma_0$$

$$t[\text{record } I_1.T_1; I_2.T_2; \dots; I_n.T_n \text{ end}]_{\varsigma_0} =$$

if distinct( $I_1, \dots, I_n$ ) then

$$\text{let } (r_1, \varsigma_1) = t[T_1]_{\varsigma_1} \text{ in }$$

$$\text{if } (\nu, \varsigma_{n+1}) = \text{newtag } \varsigma_n \text{ then}$$

$$((\nu.\text{record}:I_1.T_1, \dots, I_n.T_n), \varsigma_{n+1})$$

$$t[\uparrow I]_{\varsigma_0} = \text{if } [type:r]::[I] \text{ then}$$

$$((\nu, \uparrow r), \varsigma_2)$$

$$\text{where } (\nu, \varsigma_2) = \text{newtag } \varsigma_0$$

### 6.3. Expressions

$$e^* \in Exp^* \rightarrow U_* \rightarrow Md^*$$

$e[N]_{\varsigma} = [\text{const}:\text{int}; \text{sub}:\text{n}]$   
where  $n = \text{value}[N]$

$e[I]_{\varsigma} = \text{if } [I]::[\text{type}\ r] \text{ then } \text{error}$   
else if  $\varsigma[]::[\text{proc}:\dots]$  then  $\text{error}$   
else if  $\varsigma[]::[\text{proc}]\text{proc}$  then  $\text{error}$   
else  $\varsigma[]$

$e[O\ E]_{\varsigma} = o[O](e[E]_{\varsigma})$

$e[E_0 \Omega E_1]_{\varsigma} = w[\Omega](e[E_0]_{\varsigma})(e[E_1]_{\varsigma})$

$e[I(E_1, \dots, E_n)]_{\varsigma} =$   
if  $\varsigma[I]::[\text{a.fun}:\mu_1, \dots, \mu_n;r]$  then  
  if passable $((\mu_1, \dots, \mu_n), e^*[E_1, \dots, E_n]_{\varsigma})$  then  $[\text{val}:r]$   
  if  $\text{if passable } ((\mu_1, \dots, \mu_n), e^*[E_1, \dots, E_n]_{\varsigma})$  then  $[\text{val}:r]$   
  if  $\varsigma[I] = [\text{sfun}]$  then  $s[I](e[E_1, \dots, E_n]_{\varsigma})$

$e[E\ I]_{\varsigma} = \text{if } type(e[E])::[\nu.\text{record}(\dots, I, \dots)] \text{ then}$   
  if isval[e[E]] then  $[\text{val}:r]$  else  
    if isvar[e[E]] then  $[\text{var}:r]$

$e[E_0[E_1]]_{\varsigma} = \text{if } [\nu.\text{array } n \text{ of } r]:\text{type } e[E_0] \text{ then}$   
  if compatible $r, type\ e[E_0]_{\varsigma}$  then  
    if isvar e[E\_0]\_{\varsigma} then  $[\text{var}:r]$  else  
      if isval e[E\_0]\_{\varsigma} then  $[\text{val}:r]$

$e[\uparrow E]_{\varsigma} = \text{if } type(e[E])::[\nu, \uparrow r] \text{ then } [\text{var}:r]$

$e^* \in Exp^* \rightarrow U_* \rightarrow Md^*$

$e^*[(E_1, \dots, E_n)]_{\varsigma} = (e[E_1]_{\varsigma}, \dots, e[E_n]_{\varsigma})$

## 6.4. Statements

190

 $k \in (Blo + Decl) \rightarrow Id^*$  $g \in Stmt \rightarrow U_s \rightarrow T$ 

$$g[E_1 := E_2] \varsigma = assignable(e[E_1] \varsigma | e[E_2] \varsigma)$$

$$\begin{aligned} g[if\ E\ then\ \Gamma_1\ else\ \Gamma_2\ f] \varsigma &= isbool(type\ e[E]) \wedge \\ &\quad distinct(j[\Gamma_1]) \wedge \\ &\quad distinct(j[\Gamma_2]) \wedge \\ &\quad c[\Gamma_1] \varsigma \left[ \dots, true, \dots \right] / j[\Gamma_1] \wedge \\ &\quad c[\Gamma_2] \varsigma \left[ \dots, true, \dots \right] / j[\Gamma_2] \end{aligned}$$
 $g[dummy] \varsigma = true$ 

$$\begin{aligned} g[while\ E\ do\ \Gamma\ od] \varsigma &= \\ g[repeat\ \Gamma\ until\ E] \varsigma &= isbool(type\ e[E]) \wedge \\ &\quad distinct(j[\Gamma]) \wedge \\ &\quad c[\Gamma] \varsigma \left[ \dots, true, \dots \right] / j[\Gamma] \end{aligned}$$
 $g[goto\ N] \varsigma = \varsigma[N]$ 

$$\begin{aligned} g[I(E_1, \dots, E_n)] \varsigma &= \text{if } s[I] : proc(\mu_1, \dots, \mu_n) \text{ then} \\ &\quad passable(\mu_1, \dots, \mu_n)(e[E_1], \dots, E_n] \varsigma) \text{ else} \\ &\quad \text{if } s[I] : proc \text{ then} \\ &\quad sp[s[I](e[E_1], \dots, E_n)] \varsigma \end{aligned}$$

$$\begin{aligned} g[\Delta_1, \dots, \Delta_n, begin\ \Gamma\ end] \varsigma &= distinct(j[\Gamma]) \wedge \\ &\quad c[\Gamma] \varsigma \left[ \Delta_1, \dots, \Delta_n \right] \varsigma \\ \text{where } \varsigma &= \varsigma \left[ \dots, true, \dots \right] / j[\Gamma] \end{aligned}$$
 $j \in Com \rightarrow Num^*$ 

$$\begin{aligned} j[N:\Theta] &= N \\ j[\Gamma_1; \Gamma_2] &= j[\Gamma_1] \cdot j[\Gamma_2] \\ j[\Theta] &= \emptyset \end{aligned}$$
 $c \in Com \rightarrow U_s \rightarrow T$ 

$$\begin{aligned} c[N:\Theta] \varsigma &= g[\Theta] \varsigma \\ c[\Gamma_1; \Gamma_2] \varsigma &= c[\Gamma_1] \varsigma \wedge c[\Gamma_2] \varsigma \\ c[\Theta] \varsigma &= g[\Theta] \varsigma \end{aligned}$$

## 7. Dynamic Semantics of LS

## 7.1. Auxiliary functions

 $apply \in G \rightarrow V \rightarrow G$  $apply \gamma\epsilon = \lambda X.\gamma X\epsilon$  $change \in G_d \rightarrow (S \rightarrow S) \rightarrow G_d$ 
 $change \gamma f = \lambda X.\text{let } \theta = \gamma X \text{ in}$   
 $\quad \text{if } \gamma \in C \text{ then } \lambda \sigma.\gamma(f\sigma)$   
 $\quad \text{else } \lambda \epsilon.\text{change}(\gamma\epsilon)f$ 
 $update \in G \rightarrow G$  $update \gamma = \lambda X\alpha.\text{change}(\gamma X)(\lambda \sigma.\sigma[\epsilon/\alpha])$  $store \in (S \rightarrow G_d) \rightarrow G_d$ 
 $store f = \text{if } f \in (S \rightarrow C) \text{ then } \lambda \sigma.\text{foo}$   
 $\quad \text{else } \lambda t.\text{store}(\lambda \sigma.\text{foo})$ 
 $content \in G \rightarrow G$  $content \gamma = \lambda X\alpha.\text{store}(\lambda \sigma.(\gamma X)(\sigma\alpha))$  $cond \in G \rightarrow G \rightarrow G$  $cond\gamma_1\gamma_2 = \lambda X\epsilon. \text{if } (\epsilon \mid I) = 0 \text{ then } \gamma_1 X \text{ else } \gamma_2 X$  $binop \in Op \rightarrow G \rightarrow G$  $binop[\Omega]\gamma = \lambda X\epsilon_1.\gamma X(\psi[\Omega]\epsilon_1\epsilon_2)$  $unop \in Mop \rightarrow G \rightarrow G$  $unop[O]\gamma = \lambda X\epsilon_1.\gamma X(O[\sigma])\epsilon_1$  $Sp \in Id \rightarrow Md^* \rightarrow G \rightarrow G$  $Sp[\text{print}](\mu)\gamma = \lambda X\epsilon.\text{store}(\lambda \sigma.(\gamma X)(\sigma^{\#1}, \sigma^{\#2}, \sigma^{\#3}, \epsilon))$ 
 $Sp[\text{new}](\mu)\gamma = \lambda X\alpha. \text{if } [\text{var}.\tau]:\mu \text{ then}$   
 $\quad \text{change}(\gamma X)(\lambda \sigma. \text{let } (\dot{a}, \dot{\sigma}) = \text{heap}[\sigma] \text{ in } \dot{\sigma}[\dot{a}/a])$   
 $\text{else if } [\text{var}.\tau]:\mu \text{ then}$   
 $\quad \text{change}(\gamma X)(\lambda \sigma. \text{let } (\dot{a}, \dot{\sigma}) = \text{heap}[\sigma] \text{ in } \dot{\sigma}[\dot{a}/\alpha])$ 
 $Sf \in Id \rightarrow G \rightarrow G$  $Sf[\text{eof}]\gamma = \lambda X.\text{store}(\lambda \sigma. \text{if } |\sigma|^{\#2}| = 0 \text{ then } \gamma 0 \sigma \text{ else } \gamma 1 \sigma)$ 
 $Sf[\text{read}]\gamma = \lambda X.\text{store}(\lambda \sigma. \text{if } \sigma^{\#2}:\epsilon.w \text{ then } \gamma \epsilon(\sigma^{\#1}, w, \sigma^{\#3})$   
 $\text{else error})$

$$verifys \in I \rightarrow I \rightarrow G \rightarrow G$$

$\text{verifys } t_0 t_1 \gamma = \lambda X \epsilon. \text{ If } t_1 \leq \epsilon \wedge \epsilon \leq t_2$   
 then  $\gamma X t$   
 else adjust invalidindex(args  $\gamma X \epsilon$ )

$$\text{verifsyn} \in G \rightarrow G$$

$\text{verifsyn } \gamma = \lambda X \epsilon. \text{ If } \epsilon = \text{nil} \text{ then adjust dereferror(args } \gamma X \epsilon \text{) else } \gamma X \epsilon$

$$\text{index} \in G \rightarrow G$$

$$\text{index } \gamma = \lambda X \alpha. \gamma X(\alpha \iota)$$

$$\text{select} \in Id \rightarrow G \rightarrow G$$

$$\text{select}[I] \gamma = \lambda X \alpha. \gamma X(\alpha[I])$$

$$\text{args} \in G_d \rightarrow N$$

$\text{args } \gamma = \text{If } n = 0 \text{ then } \gamma \text{ else adjust}(\lambda \epsilon \gamma)(n - 1)$

$$\text{adjust} \in G_d \rightarrow N \rightarrow G_d$$

$\text{adjust } \gamma n = \text{If } n = 0 \text{ then } \gamma \text{ else adjust}(\lambda \epsilon \gamma)(n - 1)$

$$\text{new} \in U \rightarrow (L_r \times U)$$

$\text{new } \rho = [a, \rho[\text{true}/a]]$   
 for an  $\alpha$  such that  $\rho \alpha = \text{false}$ .

We do not further specify  $\text{new}$ ; any continuous function satisfying the above condition may be substituted.

$$\text{newl} \in Ty \rightarrow U \rightarrow (Lu \times U)$$

$$\begin{aligned} \text{newl } r \rho &= \text{if } r::[\nu_1.\text{array}[\nu_2.\text{sub}, \iota_1, \iota_2]; f] \\ &\quad \text{then newa } \iota_1 \text{ to } f \rho \\ &\quad \text{else if } r::\nu.\text{record}[\iota_1, r_1, \dots, \iota_n, r_n] \text{ then } (\alpha, \rho_n) \\ &\quad \quad \text{where } \alpha = \perp[\alpha_1/\iota_1, \dots, \alpha_n/\iota_n], \\ &\quad \quad (\alpha_i, \rho_i) = \text{newl } r_i \rho_{i-1}, \rho_0 = \rho \\ &\quad \quad \text{else new } \rho \end{aligned}$$

$$\text{newa} \in I \rightarrow I \rightarrow Ty \rightarrow U \rightarrow (Lu \times U)$$

$$\begin{aligned} \text{newa } \iota_1 \iota_2 r \rho &= \text{if } \iota_1 = \iota_2 \text{ then } \perp[\alpha/\iota_1], \rho \\ &\quad \text{where } (\alpha, \rho) = \text{newl } r \rho \\ &\quad \text{else } (\alpha_1[\alpha_2/\iota_1], \rho_2) \\ &\quad \quad \text{where } (\alpha_1, \rho_1) = \text{newa}(\iota_1 + 1) \iota_2 r \rho \\ &\quad \quad (\alpha_2, \rho_2) = \text{newl } r \rho_1 \end{aligned}$$

$$\text{cleara} \in I \rightarrow I \rightarrow Ty \rightarrow Lu \rightarrow G \rightarrow G$$

$\text{cleara } \iota_1 \iota_2 r \alpha = \text{if } \iota_1 = \iota_2 \text{ then clear } r(\alpha \iota_1) \text{ else}$   
 $\lambda \gamma. \text{clear } r(\alpha \iota_1); \text{clear}(\iota_1 + 1) \iota_2 r \gamma$

$$\text{clear} \in Ty \rightarrow Lu \rightarrow G \rightarrow G$$

$$\begin{aligned} \text{clear } r \alpha \gamma &= \text{if } r::[\nu_1.\text{array}[\nu_2.\text{sub}, \iota_1, \iota_2]; f] \\ &\quad \text{then clear } \iota_1 \iota_2 f \alpha \gamma \\ &\quad \text{else if } r::\nu.\text{record}[\iota_1, r_1, \dots, \iota_n, r_n] \\ &\quad \quad \text{then clear } r_i(\alpha[\iota_1]), \dots, \text{clear } r_n(\alpha[\iota_n]); \gamma \\ &\quad \quad \text{else (update } \gamma) \alpha \end{aligned}$$

$$\boxed{\text{heap } \sigma = (\alpha, \sigma[0/\alpha])}$$

for an  $\alpha$  such that  $\sigma\alpha = \text{unused}$ .

$$\boxed{\text{heap } \ell \in I \rightarrow Ty \rightarrow S \rightarrow (Lv_d \times S)}$$

$\text{heap} \alpha \cdot \ell_1 \cdot \ell_2 \cdot r \cdot \sigma = \text{if } \ell_1 = \ell_2 \text{ then } (\perp[\alpha/\ell_1], \sigma)$   
 where  $(\alpha, \sigma) = \text{heap} \cdot r \cdot \sigma$   
 else  $(\alpha_1[\alpha_2/\ell_1], \sigma_2)$   
 where  $(\alpha_1, \sigma_1) = \text{heap} \alpha(\ell_1 + 1) \cdot \ell_2 \cdot r \cdot \sigma$   
 $(\alpha_2, \sigma_2) = \text{heap} \cdot r \cdot \sigma_1$

$$\boxed{\text{heap} \ell \in Ty \rightarrow S \rightarrow (Lv_e \times S)}$$

## 7.2. Declarations

$\text{heap} \ell \cdot r \cdot \sigma = \text{if } r::|\nu_i:\text{array}[\nu_2:\text{sub}(\nu_1,\nu_2);f]$   
 then  $\text{heap} \alpha \cdot \ell_1 \cdot \ell_2 \cdot f \cdot \sigma$   
 else  $\text{if } r::|\nu:\text{record}[(I_1, r_1, \dots, I_n, r_n)]| \text{ then } (\alpha, \sigma_n)$   
 where  $\alpha = \perp[\alpha_1/I_1, \dots, \alpha_n/I_n]$ ,  
 $(\alpha, \sigma_i) = \text{heap} \cdot r_i \cdot \sigma_{i-1}, \sigma_0 = \sigma$   
 else  $\text{heap} \sigma$

$$\boxed{\text{enter } \ell \in B \rightarrow G \rightarrow G}$$

$$\boxed{\text{enter } n \cdot \gamma = \lambda X. \gamma(upsX \cdot n)}$$

$$\boxed{ups \in X \rightarrow B \rightarrow X}$$

$$upsX \cdot n = \text{fix } \lambda \dot{X}. (x, y, X, n)$$

where  $x = \lambda m. \text{if } m < n \text{ then } x^{\#1}m \text{ else }$

$\text{if } m = n \text{ then } \dot{X}$

where  $y = \lambda m. \text{if } m < n \text{ then } x^{\#1}m \text{ else }$

$\text{if } m = n \text{ then } \text{succ}(x^{\#2}x^{\#4})$

$$\boxed{\text{exit } \ell \in G \rightarrow G}$$

$$\boxed{\text{exit} \gamma = \lambda X. \gamma(X^{\#3})}$$

$$\boxed{\text{level } \ell \in U \rightarrow B}$$

$$\boxed{\text{level} \rho = \rho^{\#7}}$$

$$\boxed{\text{next } \ell \in U \rightarrow U}$$

$$\boxed{\text{next } \rho = (\rho^{\#1}, \dots, \rho^{\#5}, \rho^{\#6} + 1, \lambda \alpha. \text{false})}$$

$$\boxed{D \in Decl \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G}$$

$D[\text{const } \Delta_1; \dots; \Delta_n] \cdot \rho \cdot \gamma = \gamma$   
 $D[\text{type } \Delta_1; \dots; \Delta_n] \cdot \rho \cdot \gamma = \gamma$   
 $D[\text{var } \Delta_1; \dots; \Delta_n] \cdot \rho \cdot \gamma = D_v[\Delta_1; \dots; \Delta_n] \cdot \rho \cdot \gamma$   
 $D[\text{procedure } I(\Pi_1, \dots, \Pi_n); \Theta] \cdot \rho \cdot \gamma = \gamma$   
 $D[\text{function } U(\Pi_1, \dots, \Pi_n); T; \Theta] \cdot \rho \cdot \gamma = \gamma$

$$\boxed{D^* \in Decl^* \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G}$$

$$\boxed{D^*[\Delta_1; \dots; \Delta_n] \cdot \rho \cdot \gamma = D^*[\Delta_1; \dots; \Delta_{n-1}] \cdot \rho; D[\Delta_n] \cdot \rho \cdot \gamma}$$

$$\boxed{D_v \in V \text{dat} \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G}$$

$$D_v[I; T] \cdot \rho \cdot \gamma = \text{if } [\text{var } r] :: \zeta[] \text{ then } \text{clear } r(\rho[I]) \cdot \gamma$$

$$\mathcal{D}_o^* \in V^{dec^*} \rightarrow U_o \rightarrow U \rightarrow G \rightarrow G^*$$

$$\mathcal{D}_o^*[\Delta_1; \dots; \Delta_n]_o\rho\gamma = \mathcal{D}_o^*[\Delta_{o1}; \dots; \Delta_{on-1}]_o\rho; D_o[\Delta_{on}]_o\rho\gamma$$

$$\mathcal{F} \in Decl \rightarrow U_e \rightarrow U \rightarrow G \rightarrow G^*$$

$$\begin{aligned} \mathcal{F}[\text{const } \Delta_1; \dots; \Delta_n]_f\rho &= \emptyset \\ \mathcal{F}[\text{type } \Delta_1; \dots; \Delta_n]_f\rho &= \emptyset \\ \mathcal{F}[\text{var } \Delta_1; \dots; \Delta_n]_f\rho &= \emptyset \end{aligned}$$

$$\mathcal{F}[\text{procedure } I(\Pi_1, \dots, \Pi_n); \Theta]_f\rho = (\pi)$$

where  $\pi = \lambda \gamma . enter(n+1); P^*[\Pi_1, \dots, \Pi_n]_f\rho_1;$

$B[\Theta]_f\rho_1; exit \gamma$   
where  $\langle \mu_1, \dots, \mu_n \rangle, \varsigma_1 = P^*[\Pi_1, \dots, \Pi_n]_f\rho$

where  $\rho_1 = Q[\Pi_1, \dots, \Pi_n]_f\rho(\text{next } \rho)$   
where  $n = level \rho$

$$\mathcal{F}[\text{function } I_i(\Pi_1, \dots, \Pi_n); f_i, \Theta]_f\rho = (\pi)$$

where  $\pi = \lambda \gamma . enter(n+1); P^*[\Pi_1, \dots, \Pi_n]_f\rho_1;$   
 $B[\Theta]_f\rho_1; return[I_i]_f\rho_2; exit \gamma$   
where  $\varsigma_2 = \varsigma_1[\text{afun}; \mu_1, \dots, \mu_n, \tau]/I_i$

where  $\langle \mu_1, \dots, \mu_n \rangle, \varsigma_1 = P^*[\Pi_1, \dots, \Pi_n]_f\rho$   
where  $\rho_1 = Q[\Pi_1, \dots, \Pi_n]_f\rho(\text{next } \rho)$   
where  $n = level \rho$

$$P^*[\Pi_1, \dots, \Pi_n]_f\rho\gamma = P[\Pi_n]_f\rho; P[\Pi_{n-1}]_f\rho; \dots; P[\Pi_1]_f\rho\gamma$$

$$Q \in Par \rightarrow U_e \rightarrow U \rightarrow U$$

$$\mathcal{F}^* \in Decl^* \rightarrow U_e \rightarrow U \rightarrow G^*$$

$$\mathcal{F}^*[\Delta_1; \dots; \Delta_n]_f\rho = \mathcal{F}[\Delta_1]_f\rho; \dots; \mathcal{F}[\Delta_n]_f\rho$$

$$\begin{aligned} Q[\text{var } I_1; I_2]_f\rho &= \\ Q[I_1; I_2]_f\rho &= \delta[a/I_1] \quad \text{where } (a, \rho) = new \rho \end{aligned}$$

$$\mathcal{V} \in Decl \rightarrow U_e \rightarrow U \rightarrow U$$

$$\begin{aligned} \mathcal{V}[\text{const } \Delta_1; \dots; \Delta_n]_v\rho &= \rho \\ \mathcal{V}[\text{type } \Delta_1; \dots; \Delta_n]_v\rho &= \rho \\ \mathcal{V}[\text{var } \Delta_1; \dots; \Delta_n]_v\rho &= \mathcal{V}_v[\Delta_{v1}; \dots; \Delta_{vn}]_v\rho \\ \mathcal{V}[\text{procedure } I(\Pi_1, \dots, \Pi_n); \Theta]_v\rho &= \rho \\ \mathcal{V}[\text{function } I(\Pi_1, \dots, \Pi_n); \Theta]_v\rho &= \rho \end{aligned}$$

$$\mathcal{V}^* \in Decl^* \rightarrow U_e \rightarrow U \rightarrow U$$

$$\mathcal{V}^*[\Delta_1; \dots; \Delta_n]_v\rho = \mathcal{V}^*[\Delta_2; \dots; \Delta_n]_v\rho; \mathcal{V}[\Delta_1]_v\rho$$

$$\mathcal{V}_v \in V^{dec} \rightarrow U_e \rightarrow U \rightarrow U$$

$$\begin{aligned} \mathcal{V}_v[I; \mathcal{X}]_v\rho &= \text{if } [var; \tau]; \mathcal{X}[I] \text{ then } \delta[a/I] \\ &\quad \text{where } (a, \rho) = new \tau \rho \end{aligned}$$

$$P \in Par \rightarrow U_e \rightarrow U \rightarrow G \rightarrow G$$

$$\begin{aligned} P[I_1; I_2]_f\rho\gamma &= \lambda c . (\text{update } \gamma)(c(\rho[I_1])) \\ P[\text{var } I_1; I_2]_f\rho\gamma &= \lambda c . (\text{update } \gamma)(c(\rho[I_1])) \end{aligned}$$

$$P^*[\Pi_1, \dots, \Pi_n]_f\rho\gamma = P[\Pi_n]_f\rho; P[\Pi_{n-1}]_f\rho; \dots; P[\Pi_1]_f\rho\gamma$$

$$\mathcal{Q}^* \in \text{Par}^* \rightarrow U_s \rightarrow U \rightarrow U$$

$$\mathcal{Q}^*(\Pi_1, \dots, \Pi_n) \cdot \rho = \mathcal{Q}^*(\Pi_1, \dots, \Pi_{n-1}) \cdot (\mathcal{Q}(\Pi_n) \cdot \rho)$$

## 7.3. Expressions

$$\mathcal{E}, \hat{\mathcal{E}} \in \text{Exp} \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G$$

$$\mathcal{E}[E] \cdot \rho \cdot \gamma = \text{if } e[E] \cdot \cdot \cdot [\text{const}; \cdot; \cdot] \text{ then apply } \gamma \cdot \epsilon \text{ else } \hat{\mathcal{E}}[E] \cdot \rho \cdot \gamma$$

$$\mathcal{E}[N] \cdot \rho \cdot \gamma = \text{apply } \gamma(\text{value}[N])$$

$$\begin{aligned} \mathcal{E}[I] \cdot \rho \cdot \gamma &= \text{if } s[I] = [\text{sfun}] \text{ then } S[I] \cdot \gamma \text{ else} \\ &\quad \text{if } s[I] : [\text{afun}; () \cdot \text{tau}] \text{ then } (\rho[I])^\# \cdot \gamma \text{ else} \\ &\quad \text{if } s[I] : [\text{pfun}; () \cdot \text{tau}] \text{ then } (\rho[I])^\# \cdot \gamma \text{ else apply } \gamma(\rho[I]) \end{aligned}$$

$$\hat{\mathcal{E}}[\text{O } E] \cdot \rho \cdot \gamma = R[E] \cdot \rho; \text{unop}[\text{O}] \cdot \gamma$$

$$\hat{\mathcal{E}}[E_0 \Omega E_1] \cdot \rho \cdot \gamma = R[E_0] \cdot \rho; R[E_1] \cdot \rho; \text{binop}[\Omega] \cdot \gamma$$

$$\hat{\mathcal{E}}[E \dagger E_1] \cdot \rho \cdot \gamma = R[E] \cdot \rho; \text{verify}[\gamma]$$

$$\hat{\mathcal{E}}[I(E_1, \dots, E_n)] \cdot \rho \cdot \gamma = \text{if } s[I] = [\text{sfun}] \text{ then } \mathcal{E}'[E_1, \dots, E_n] \cdot \rho; Sf[I] \cdot \gamma \text{ else}$$

$$\text{if } s[I] : [\text{afun}; \mu_1 \dots \mu_n; \cdot] \text{ then } A[I] \cdot \rho; \mu_1 \dots \mu_n \cdot \rho; (\rho[I])^\# \cdot \gamma \text{ else}$$

$$\text{if } s[I] : [\text{pfun}; \mu_1 \dots \mu_n; \cdot] \text{ then } A[I] \cdot \rho; \mu_1 \dots \mu_n \cdot \rho; (\rho[I])^\# \cdot \gamma$$

$$\hat{\mathcal{E}}[E_0[E_1]] \cdot \rho \cdot \gamma = \mathcal{E}[E_0] \cdot \rho; A[E_1] \cdot [\text{val}; \text{to}; \rho; \text{index} \cdot \gamma]$$

$$\text{where type}[\epsilon(E_0)] \cdot [\nu; \text{array}; \text{to}; \gamma]$$

$$\hat{\mathcal{E}}[E.I] \cdot \rho \cdot \gamma = \mathcal{E}[E] \cdot \rho; \text{select}[I] \cdot \gamma$$

$$\mathcal{E}^* \in \text{Exp}^* \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G$$

$$\mathcal{E}^*[E_1, \dots, E_n] \cdot \rho \cdot \gamma = \mathcal{E}[E_1] \cdot \rho, \mathcal{E}[E_2] \cdot \rho, \dots, \mathcal{E}[E_n] \cdot \rho \cdot \gamma$$

$$\mathcal{A} \in \text{Exp} \rightarrow U_s \rightarrow M_d \rightarrow U \rightarrow G \rightarrow G$$

$$\mathcal{A}^*[E_1, \dots, E_n] \cdot \rho \cdot \gamma = \mathcal{A}[E_1] \cdot \rho, \mathcal{A}[E_2] \cdot \rho, \dots, \mathcal{A}[E_n] \cdot \rho \cdot \gamma$$

## 7.4. Statements

$$\begin{aligned}
 C \in Com &\rightarrow U_s \rightarrow U \rightarrow G \rightarrow C \\
 C[N:\Theta]\varsigma\rho\gamma &= B[\Theta]\varsigma\rho\gamma \\
 C[\Theta]\varsigma\rho\gamma &= \beta[\Theta]\varsigma\rho\gamma \\
 C[\Gamma_0;\Gamma_1]\varsigma\rho\gamma &= C[\Gamma_0]\varsigma\rho\gamma; C[\Gamma_1]\varsigma\rho\gamma
 \end{aligned}$$

$$J \in Com \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G'$$

$$\begin{aligned}
 C[\Gamma]\varsigma\rho\gamma &= C[\Gamma]\hat{\beta}\rho\gamma \\
 \text{where } \hat{\beta} &= \varsigma([\dots, true, \dots]) / J[\Gamma], \\
 J[N:\Theta]\varsigma\rho\gamma &= (B[\Theta]\varsigma\rho\gamma) \\
 J[\Gamma_0;\Gamma_1]\varsigma\rho\gamma &= J[\Gamma_0]\varsigma\rho\gamma; (C[\Gamma_2]\varsigma\rho\gamma) \\
 J[\Gamma_2]\varsigma\rho\gamma &= J[\Gamma_2]\varsigma\rho\gamma
 \end{aligned}$$

$$J \in Com \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G'$$

$$\begin{aligned}
 J[N:\Theta]\varsigma\rho\gamma &= (B[\Theta]\varsigma\rho\gamma) \\
 J[\Gamma_0;\Gamma_1]\varsigma\rho\gamma &= J[\Gamma_0]\varsigma\rho\gamma; (C[\Gamma_2]\varsigma\rho\gamma) \\
 J[\Gamma_2]\varsigma\rho\gamma &= J[\Gamma_2]\varsigma\rho\gamma
 \end{aligned}$$

$$\theta \in Stmt \rightarrow U_s \rightarrow U \rightarrow G \rightarrow G$$

$$\begin{aligned}
 B[E_0 := E_1]\varsigma\rho\gamma &= L[E_0]\varsigma\rho; A[E_1]\varsigma[{\it val}; r]\rho; update\ \gamma \\
 \text{where } e[E_0]\varsigma &::: [var : r] \\
 B[if\ E\ then\ F\ o\ else\ G]\varsigma\rho\gamma &= R[E]\varsigma\rho; Cond(C[\Gamma_0]\varsigma\rho\gamma, C[\Gamma_1]\varsigma\rho\gamma) \\
 B[while\ E\ do\ F]\varsigma\rho\gamma &= fix(\lambda\ t. R[E]\varsigma\rho; cond(C[\Gamma]\varsigma\rho\gamma, \gamma)) \\
 B[repeat\ F\ until\ E]\varsigma\rho\gamma &= fix(\lambda\ t. C[\Gamma]\varsigma\rho; R[E]\varsigma\rho; cond(r, t)) \\
 B[goto\ N]\varsigma\rho\gamma &= adjust(\rho[N])(args\ \gamma - args(\rho[N]))
 \end{aligned}$$

**1. Abstract syntax**

$n \in N$	Numerals
$l \in L$	Labels
$p \in P = I_1 \dots I_n$	Programs
$Cde = N + P + I$	Machine Code
$t \in I =$	
$LIT\ n \mid LOAD\   ADDR\ n\ n \mid STORE\  $	
$CALL\ l\ n_1\ n_2 \mid EXIT\mid NEW\mid$	
$NCOND\ l\   HOP\ l\   JUMP\ l\ n\mid$	
$LABSET\ l\   STOP\   UNOP\ n\mid BINOP\ n\mid$	
$EOF\mid OUTP\mid INPT$	Instructions

**2. Semantic Domains**

$x \in N$	Integers
$\sigma \subseteq S = V^*$	Stacks
$\rho \in U = L \rightarrow G$	Environments
$\delta \in D = N \rightarrow D \times N \rightarrow N \times D \times N \times G \times N$	Displays
$\mu \in M = (N \rightarrow N) \times N^* \times N^* \times N$	Memory
$\gamma \in G = D \rightarrow S \rightarrow M \rightarrow A$	Continuations
$A = N^* + \{eoferror, nostop, \dots\}$	Answers

**3. Auxiliary Functions**

$P[l_1 \dots l_n] \text{ in file} = \text{if } \text{distinct}([l_1 \dots l_n]) \text{ then } \kappa^*[l_1 \dots l_n] \rho \tau_0 \chi_0 \sigma_0 \mu_0$   
 where  $\rho_1 = \text{fix}(\lambda \rho. \rho_0[l_1 \dots l_n] \rho \gamma_0 / [l_1 \dots l_n])$

$$\boxed{up_T \in D \rightarrow G \rightarrow N \rightarrow N \rightarrow D}$$

$$\begin{aligned} up_T \delta \gamma n m i &= \text{fix}(\lambda \delta. f_i g, \delta, n, \gamma, m)) \\ &\quad \text{where } f = \lambda \nu. \text{if } \nu < n \text{ then } \delta^{#1} \nu \text{ else} \\ &\quad \quad \text{If } \nu = n \text{ then } \delta \text{ else } \frac{1}{\delta} \\ &\quad \quad \text{where } g = \lambda \nu. \text{if } \nu < n \text{ then } \delta^{#2} \nu \text{ else} \\ &\quad \quad \quad \text{If } \nu = n \text{ then } k \text{ else } \perp \\ &\quad \quad \quad \text{where } k = \delta^{#2} \delta^{#4} + i \end{aligned}$$

$$\boxed{\text{length} \in S \rightarrow N}$$

$$\text{length } \sigma = \text{if } \sigma :: x \dot{\sigma} \text{ then } 1 + \text{length } \dot{\sigma} \text{ else } 0$$

$$\boxed{down \in S \rightarrow N \rightarrow S}$$

$$\text{down } \sigma n = \text{if } \text{length } \sigma = n \text{ then } \sigma \text{ else } \text{down}(\text{pop } \sigma)$$

$$\boxed{\text{pop} \in S \rightarrow S}$$

$$\text{pop } \sigma = \text{if } \sigma :: x \dot{\sigma} \text{ then } \dot{\sigma}$$

**4. Semantic Equations**

$$\boxed{P \in P \dashv N^* \dashv (N^* + E_i)}$$

For a program the initial environment, continuation, display, stack, and memory are given as:

$$\rho_0 = \perp$$

$$\gamma_0 = \text{fix } (\lambda \delta. \langle \perp[\delta/0], \perp[0/0], \perp, 0, \perp, 0 \rangle)$$

$$\sigma_0 = \langle \rangle$$

$$\mu_0 = \langle \perp, \text{infile}, \langle \rangle, 0 \rangle$$

$$L \in P \rightarrow U \rightarrow G \rightarrow G^*$$

$$L[I_1, \dots, I_n] \rho\gamma = \text{if } I_1 : [\text{LABSET } l] \\ \text{then } \langle M[V_2, \dots, I_n] \rho\gamma \rangle \cdot L[I_2, \dots, I_n] \rho\gamma \\ \text{else } L[I_2, \dots, I_n] \rho\gamma$$

$$I \in P \rightarrow N^*$$

$$\ell[I_1, \dots, I_n] = \text{if } I_1 : [\text{LABSET } l] \text{ then } \langle n \rangle \ell[I_2, \dots, I_n] \\ \text{else } \ell[I_2, \dots, I_n]$$

$$M \in I \rightarrow U \rightarrow G \rightarrow G$$

$$M[LIT n] \rho\gamma = \lambda \delta \sigma \cdot \gamma \delta(n \sigma)$$

$$M[LOAD] \rho\gamma = \lambda \delta \sigma \mu. \text{if } \sigma : x \cdot \delta \text{ then } \gamma \delta((\mu x) \cdot \delta) \mu$$

$$M[ADDR z n] \rho\gamma = \lambda \delta \sigma \mu \cdot \gamma \delta(\alpha \cdot \sigma) \mu \text{ where } \alpha = (\delta^{*2} n) + z$$

$$M[STORE] \rho\gamma = \lambda \delta \sigma \mu. \text{if } \sigma : x \cdot y \cdot \delta \text{ then } \gamma \delta \sigma(\mu^* | x / y)$$

$$M[CALL l n m k] \rho\gamma = \lambda \delta \sigma \mu [l] \cdot (\text{upr } \delta \cdot \gamma \cdot n \cdot (\text{length } \sigma - m)k) \sigma$$

$$M[EXIT] \rho\gamma = \lambda \delta \cdot \gamma(\delta^{*3}) \quad \text{where } \gamma = \delta^{*5}$$

$M[NCOND] \rho\gamma = \lambda \delta \sigma. \text{if } \sigma :: x \cdot \delta \text{ then } (\text{if } x \neq 0 \text{ then } \rho[l] \delta \sigma \text{ else } \gamma \delta \delta)$

$$M[HOP l] \rho\gamma = \rho[l]$$

$$M[JUMP l n] \rho\gamma = \lambda \delta \sigma \cdot \rho[z](\delta^{*1} n)(\text{down } \sigma(\delta^{*1} n)^{*6})$$

$$M[LABSET l] \rho\gamma = \gamma$$

$$M[STOP] \rho\gamma = \lambda \delta \sigma \mu \cdot \mu^{*3}$$

$$M[NEW] \rho\gamma = \lambda \delta \sigma \mu \cdot \gamma \delta(\mu^{*4} \sigma)(\mu^{*1}, \mu^{*2}, \mu^{*3}, \mu^{*4} + 1)$$

$$M[U NOP n] \rho\gamma = \lambda \delta \sigma. \text{if } \sigma :: z \cdot \delta \text{ then } \gamma \delta((O[n]z) \cdot \delta)$$

$$M[BINOP n] \rho\gamma = \lambda \delta \sigma. \text{if } \sigma :: z \cdot y \cdot \delta \text{ then } \gamma \delta((\psi[n]y) \cdot z) \cdot \delta$$

$$M[EOF] \rho\gamma = \lambda \delta \sigma \mu. \text{if } \mu^{*2} = \langle \rangle \text{ then } \gamma \delta(0 \cdot \sigma) \mu \text{ else } \gamma \delta(1 \cdot \sigma) \mu$$

$$M[OUTP] \rho\gamma = \lambda \delta \sigma \mu. \text{if } \sigma :: z \cdot \delta \text{ then } \gamma \delta \delta / \mu$$

$$\text{where } \mu = \langle \mu^{*1}, \mu^{*2}, \mu^{*3} \cdot z \rangle$$

$$M[INPT] \rho\gamma = \lambda \delta \sigma \mu. \text{if } \mu^{*2} :: x \cdot \mu \text{ then }$$

$$\gamma \delta(x \cdot \sigma) \langle \mu^{*1}, \mu, \mu^{*3} \rangle \text{ else coerror}$$

$O$  and  $\psi$  are defined as in LS.

### Appendix 3. The Scanner

#### 1. Logical basis

##### 1.1. Definition of the miero syntax

```

rulefile (scanner)
constant      color, period, equal, null_sequence,
idtoken, errortoken,
periodtoken, dperiodtoken, numbertoken, length, value,
void, syn, sem, tops, names, values,
decomettoken, colontoken;
from true infer
colon ≠ period ∧ colon ≠ equal ∧ period ≠ equal ∧
length ≠ value ∧ tops ≠ names ∧ syn ≠ sem ∧
sem ≠ void ∧ syn ≠ void;

%definition of the scanning function %

psit: replace PSI([tab, s) where eos(s) by null_sequence;
psit2: replace PSI([tab, s) where letter(hd(s)) by PHIdent(tab, list(hd(s)), t([s]));
psit3: replace PSI([tab, s) where digit(hd(s)) by PHInumber(tab, list(hd(s)), t([s]));
psit4: replace PSI([tab, s) where delim(hd(s)) by PHIdelim(tab, list(hd(s)), t([s]));
psit5: replace PSI([tab, s) where colon = hd(s) = hd(s) by PHIconc([tab, list(hd(s)), t([s));
psit6: replace PSI([tab, s) where period = hd(s) by PHIPer(tab, list(hd(s)), t([s));
psit7: replace PSI([tab, s) where ~letter(hd(s))

%scanning identifiers, numbers and delimiters %
phi1: replace PHIIdent(tab, s1, s2) where letter(hd(s2))
      by PHIdent(tab, concat(s1, list(hd(s2))), t([s2));
phi2: replace PHIIdent(tab, s1, s2) where digit(hd(s2))
      by PHIdent(tab, concat(s1, list(hd(s2))), t([s2));
phi3: replace PHIIdent(tab, s1, s2) where ~digit(hd(s2))
      by PHIIdent(tab, concat(list(SIdent(tab, s1)), PSI(STIdent(tab, s1), s2)));
phi4: replace PHIIdent(tab, s1, null_sequence)
      by list(SIdent(tab, s1));
phi5: replace PHICol(tab, list(color), null_sequence) by list(SCol(1));
phi6: replace PHICol(tab, list(colon), s) where hd(s) = equal
      by concat(list(SCol(2)), PSI([tab, t([s))));
phi7: replace PHICol(tab, list(colon), s) where hd(s) ≠ equal
      by concat(list(SCol(1)), PSI([tab, s)));
phi8: replace PHIper(tab, list(period), null_sequence) by list(Sper(1));
phi9: replace PHIper(tab, list(period), s) where hd(s) = period
      by concat(list(Sper(2)), PSI([tab, t([s))));
phi10: replace PHIper(tab, list(period), s) where hd(s) ≠ period
      by concat(list(Sper(1)), PSI([tab, s)));
phi11: replace PHIper(tab, list(c), s) by concat(list(Sdelim([tab, c])), PSI([tab, s]));
phi12: replace PIIInumber(tab, s1, null_sequence) by list(Snumber(tab, s1));
phi13: replace PIIInumber(tab, s1, s2) where digit(hd(s2))
      by PIIInumber(tab, concat(s1, list(hd(s2))), t([s2]));

```

**phm3:** replace  $P/H/\text{Inumber}(\text{tab}, s1, s2)$  where  $\neg\text{digit}(\text{hd}(s2))$   
by  $\text{concat}(\text{list}(\text{Snumber}(\text{tab}, s1)), P/S_I(\text{tab}, s2))$ ;

%semantic routines attached to tokens %

**Scl:** replace  $S\text{col}(1)$  by  $\text{mkvoidtoken}(\text{colon token})$ ;

**Scl2:** replace  $S\text{col}(2)$  by  $\text{mkvoidtoken}(\text{becomes token})$ ;

**Sp1:** replace  $S\text{per}(1)$  by  $\text{mkvoidtoken}(\text{period token})$ ;

**Sp2:** replace  $S\text{per}(2)$  by  $\text{mkvoidtoken}(\text{period token})$ ;

**Sn1:** replace  $S\text{number}(\text{tab}, s)$  by  $\text{mktoken}(\text{number token}, N\text{acr}(s))$ ;

**Sn2:** replace  $S\text{ident}(\text{tab}, s)$  where  $\kappa\text{appa}(s).\text{syn} = \text{errortoken}$   
by  $\text{mktoken}(\text{id token}, \text{apply}(\text{henter}(\text{tab}, s), s))$ ;

**Stid2:** replace  $S\text{ident}(\text{tab}, s)$  where  $\kappa\text{appa}(s).\text{syn} \neq \text{errortoken}$  by  $\kappa\text{appa}(s)$ ;

**STid1:** replace  $S\text{Tident}(\text{tab}, s)$  where  
 $\kappa\text{appa}(s).\text{syn} = \text{errortoken}$  by  $\text{henter}(\text{tab}, s)$ ;

**STid4:** replace  $S\text{Tident}(\text{tab}, s)$  where  $\kappa\text{appa}(s).\text{syn} \neq \text{errortoken}$  by  $\text{tab}$ ;

**Nscr1:** replace  $N\text{scr}(\text{hd}(c))$  by  $\text{val}(c)$ ;

**Nscr2:** replace  $N\text{scr}(\text{concat}(\text{list}(c)))$  by  $10^*N\text{scr}(c) + \text{val}(c)$ ;

**1.2. Representation functions**

**hen1:** replace  $\text{henter}(\text{tab}, s)$  where  $0 = \text{apply}(\text{tab}, s)$  by  $\text{redefine}(\text{tab}, s)$ ;

**hen2:** replace  $\text{henter}(\text{tab}, s)$  where  $0 \neq \text{apply}(\text{tab}, s)$  by  $\text{tab}$ ;

%definition of token representation %

**voidt1a:** infer  $z = \text{mkvoidtoken}(a)$  from  $z.\text{void} \wedge z.\text{syn} = a$ ;

**voidt1:** whenever  $z.\text{void}$  infer  $z = \text{mkvoidtoken}(z.\text{syn})$ ;

**mtok1a:** infer  $z = \text{mtok1}(a, b)$  from  $\neg z.\text{void} \wedge z.\text{syn} = a \wedge z.\text{sem} = b$ ;

**mtok1:** whenever  $z.\text{void}$  from  $\neg z.\text{void}$  infer  $z = \text{mtok1}(z.\text{syn}, z.\text{sem})$ ;

%identifier table %

**it1:** replace  $\text{apply}(\text{tabrep}(t), s2)$  by cases  
 $t.\text{tops} = 0 \rightarrow 0$ ;  
 $t.\text{tops} \neq 0 \wedge \text{stringrep}(\text{t.names}[t.\text{tops}]) = s2 \rightarrow t.\text{type};$   
 $t.\text{tops} \neq 0 \wedge \text{stringrep}(\text{t.names}[t.\text{tops}]) \neq s2 \rightarrow$   
 $\text{apply}(\text{tabrep}(< t, \text{tops}, t.\text{tops} - 1 >), s2) \text{ end };$

**it2:** replace  $\text{newvalue}(\text{tabrep}(t))$  by  $t.\text{tops} + 1$ ;

**it3:** replace  $\text{visible}(\text{tab})$  by  $\text{tab}.\text{tops} \geq 0$ ;

**it4:** infer  $\text{apply}([t1, s] = \text{apply}([t2, s]$  from  $t1.\text{tops} \leq t2.\text{tops} \wedge \text{apply}(t1, s) \neq 0$ );

**it5:** replace  $\text{redesign}(\text{tabrep}(t), \text{stringrep}(s))$  by  
 $\text{tabrep}(< t, \text{names}[t.\text{tops} + 1], s >, \text{tops}, t.\text{tops} + 1 >)$ ;

**1.3. Sequences**

%various implementation details %

**str1:** infer  $\text{stringrep}(z) = \text{null\_sequence}$  from  $z.\text{length} = 0$ ;

**str1a:** replace  $\text{stringrep}(< z, \text{length}, 0 >)$  by  $\text{null\_sequence}$ ;

**str2:** replace  $\text{concat}(\text{stringrep}(z), \text{list}(c))$  by  
 $\text{stringrep}(< < z, \text{length}, \text{z.length} + 1 >, \text{value},$   
 $< z.\text{value}, [\text{z.length} + 1], c >>);$

**str3:** from  $\text{stringrep}(z) = \text{stringrep}(y) \wedge$   
 $y.\text{value}[\text{z.length} + 1] = y.\text{value}[y.\text{length} + 1]$   
infer  $\text{stringrep}(< z, \text{length}, \text{z.length} + 1 >) =$   
 $\text{stringrep}(< y, \text{length}, y.\text{length} + 1 >);$

**str4:** infer  $\text{stringrep}(z) \neq \text{stringrep}(y)$  from  $\neg z.\text{length} = y.\text{length}$ ;

**str5:** from  $\text{stringrep}(< z, \text{length}, i >) \neq \text{stringrep}(< y, \text{length}, i >)$   $\wedge$   
 $i < z.\text{length}$

```

infer stringrep(z) ≠ stringrep(y);

str6: infer stringrep(z) ≠ stringrep(y) from
      z.value[z.length] ≠ y.value[y.length];
      tokenkind == periodtoken;
      errortoken, dperiodtoken, colontoken, periodtoken,
      carray == array [1.. stringlength] of char;
      string == record length: integer; value: array end;
      token ... record syn: tokenkind; sem: integer; void: boolean end;
      tokenfile == file of token;
      charfile == file of char;
      itype == record tops: integer;
      names: array [1.. itabsize] of string;
      values: array [1.. itabsize] of integer end;

var cvalid : boolean;
infile, infile0: charfile;
outfile: tokensfile;
c: char;
itab, itab0: itype;

procedure error;
entry true;
exit false;
% This procedure will print an error message. Any further action of the
program will not be covered by the proof; This allows to conclude false
as an exit condition. %
external;

function letter(c: char): boolean;
%letter(c) iff c ∈ {A, B, ..., x, y, z} no proof for this %
entry true; exit true; external;

function digit(c: char): boolean;
%digit(c) iff c ∈ {0, 1, ..., 9} %
entry true; exit true; external;

function delim(c: char): boolean;
%delim(c) iff c ∈ {;, &, V, |} =
entry true; exit true; external;

%append a character to a string %
procedure append(var s: string; c: char);
initial s == #0;
entry true;

2. The program

pascal

const stringlength = 20;
itabsize = 200;
errortoken = 0;
period == '.';
colon == ':';
equal == '=';
```

```

exit stringrep(s) = concat(stringrep(s0), istt(c));
begin
  if s.length = stringlength then error
  else begin
    s.length ← s.length + 1;
    s.value[s.length] ← c;
  end;
end; %append %

function stringequal(s1, s2: string): boolean;
entry true; exit stringequal = (stringrep(s1) = stringrep(s2));
var i: integer;
eq: boolean;

begin
  if s1.length ≠ s2.length then stringequal ← false else
  begin
    i ← 0;
    eq ← true;
    invariant
      ((eq ∧ stringrep(< s1, length, i >) = stringrep(< s2, length, i >))
       ∨ (¬eq ∧ i ≤ s1.length ∧ i > stringrep(< s2, length, i >)))
       ∧ 0 ≤ i ≤ s1.length ∧ i ≤ s2.length
       ∧ s2.length = s1.length
    while (i < (s1.length)) and eq do
    begin
      i ← i + 1;
      if s1.value[i] ≠ s2.value[i] then eq ← false;
    end;
    stringequal ← eq;
  end;
end; %stringequal %

%make a token with void semantic information %
function voidtoken(t: tokenkind): token;
entry true; exit voidtoken = mkevoidtoken(t);
var tok: token;
begin
  tok.syn ← t;
  tok.sem ← 0;

```

```

tok_void ← true;
voidtoken ← tok;
end;

%make a token with semantic information s, nonvoid%
function mktoken(t, s: integer): token;
entry true;
exit mktoken = mktoken(t, s);
var tok: token;
begin
  tok.syn ← t;
  tok.sem ← s;
  tok_void ← false;
  mktoken ← tok;
end;

procedure icenter(var tab: itype; var r: integer; s: string);
initial tab = tab0;
entry isitable(tab);
exit isitable(tab) ∧
tabrep(tab) = henter(tabrep(tab0), stringrep(s)) ∧
var i : integer;
found: boolean;
begin
  found ← false;
  i ← 0;
  invariant
    ((¬found ∧ apply(tabrep(< tab, top, i >), stringrep(s)) = errrep) ∨
     (found ∧ apply(tabrep(< tab, top, i >), stringrep(s)) = i ∧ r = i ∧ r > 0))
    ∧ i ≤ tab.top ∧ 0 ≤ i;
  while (i < (tab.top)) and not found do
  begin
    i ← i + 1;
    if stringequal(i, tab.names[i]) then
    begin
      r ← i;
      found ← true;
    end;
  end;
  if not found then

```

```

begin
  tab.topo ← tab.topo + 1;
  tab.names[tab.topo] ← s;
  r ← tab.topo;
end;
%enter
begin
  tenter(itab, t, sem, s),
  t_void ← false,
end;
%mapident %
begin %body of scanident %
  s.length ← 0;
repeat
  append(s, c);
  if eof(infile) then valid ← false else read(infile, c)
  until ((not valid) or (not(letter(c) or digit(c))))
  invariant ((valid ∧
    PHiIdent(tabrep(itab0), list(c0), infile0) =  

    V(~valid ∧
      PHiIdent(tabrep(itab0), stringrep(s), concat(list(c), infile0)))  

    concat(outfile, PS(tabrep(itab), concat(list(c), infile))) =  

    concat(outfile, PHiIdent(tabrep(itab0), list(c0), infile0))) ∧  

    V(~valid ∧ outfile =  

      concat(outfile, PHiIdent(tabrep(itab0), list(c0), infile0))) ∧  

    (valid ∨ eof(infile));
  mapident(tok, s);
  write(outfile, tok);
end; %scanident %

%scanning identifiers %
procedure scanident;
global ( var valid, c, infile, outfile, itab );
initial infile = infile0, outfile = outfile0, c = c0, itab = itab0;
entry letter(c) ∧ valid ∧ isitable(itab);
exit  ((valid ∧
  concat(outfile, PS(tabrep(itab), concat(list(c), infile))) =  

  concat(outfile, PHiIdent(tabrep(itab0), list(c0), infile0))) ∧  

  (valid ∨ eof(infile)) ∧ isitable(itab));
var s : string;
tok: token;
begin
  function keywords(s: string): token;
  entry  true;
  exit   keywords = kappa(stringrep(s));
  var l, m, r: integer;
  found: boolean;
  result: token;
  external ;
procedure mapident( var t: token; s: string);
  entry  true;
  global ( var itab );
  initial itab = itab0;
  exit   t = SIdent(tabrep(itab0), stringrep(s)) ∧
    tabrep(itab) = SIdent(tabrep(itab0), stringrep(s)) ∧
    isitable(itab);
var i : integer;
begin
  t ← keywords(s);
  if errorToken = t.syn then
begin
  t.syn ← idtoken;

```

```

  tenter(itab, t, sem, s),
  t_void ← false,
end;
%mapident %
begin %body of scanident %
  s.length ← 0;
repeat
  append(s, c);
  if eof(infile) then valid ← false else read(infile, c)
  until ((not valid) or (not(letter(c) or digit(c))))
  invariant ((valid ∧
    PHiIdent(tabrep(itab0), list(c0), infile0) =  

    V(~valid ∧
      PHiIdent(tabrep(itab0), stringrep(s), concat(list(c), infile0)))  

    concat(outfile, PS(tabrep(itab), concat(list(c), infile))) =  

    list(SIdent(tabrep(itab0), stringrep(s))) )) ∧  

    (valid ∨ eof(infile));
  mapident(tok, s);
  write(outfile, tok);
end; %scanident %

%scanning numbers %
procedure scannumber;
global ( var valid, c, infile, outfile, itab );
initial infile = infile0, outfile = outfile0, c = c0;
entry digit(c) ∧ valid;
exit  ((valid ∧
  concat(outfile, PSL(tabrep(itab), concat(list(c), infile))) =  

  concat(outfile, PHInumber(tabrep(itab), list(c0), infile0))) ∧  

  (valid ∨ eof(infile)));
var s : string;
  v: integer;
begin
  function val(c: char): integer;
  %this is ord(c) - ord('0'); take it for granted.. %
  entry digit(c); exit true; external ;
begin %body of scannumber %
  s.length ← 0;

```

```

 $v \leftarrow 0;$ 
 $\text{repeat}$ 
 $\quad \text{append}(s, c);$ 
 $\quad v \leftarrow 10v + val(c);$ 
 $\text{if } eos(\text{infile}) \text{ then } valid \leftarrow \text{false} \text{ else read}(\text{infile}, c)$ 
 $\text{until } not valid \text{ or not digit}(c)$ 
 $\text{invariant } ((valid \wedge$ 
 $\quad PHInumber(\text{tabrep}(\text{itab}), list(c0), infile) =$ 
 $\quad PHInumber(\text{tabrep}(\text{itab}), stringrep(s), concat(list(c), infile))) \wedge$ 
 $\quad (\neg valid \wedge$ 
 $\quad PHInumber(\text{tabrep}(\text{itab}), list(c0), infile) =$ 
 $\quad list(Snumber(\text{tabrep}(\text{itab}), stringrep(s)))))) \wedge$ 
 $\quad (valid \vee eos(infile)) \wedge$ 
 $\quad Nscr(stringrep(s)) = v;$ 
 $\quad write(outfile, mtoken(numberToken, v));$ 
 $\text{end : \%scannumber \%}$ 

%scanning delimiters %
procedure scolon;
global ( var valid, c, infile, outfile; itab);
initial infile == infile0, outfile == outfile0, c == c0;
entry deftm(c)  $\wedge$  cvaid;
exit  ((valid  $\wedge$ 
concat(outfile, PSI(tabrep(itab), concat(list(c), infile))) ==
concat(outfile, PHIdelim(tabrep(itab), list(c0), infile0)))  $\wedge$ 
(V(\neg cvaid  $\wedge$  outfile ==
concat(outfile0, PHIdelim(tabrep(itab), list(c0), infile0)))  $\wedge$ 
(valid  $\vee$  eos(infile)));
(outfile, mtoken(scolon, c));
begin
write(outfile, mapdelim(itab, c));
if eos(infile) then valid  $\leftarrow$  false else read(infile, c);
end ; %scandelim %

procedure scancer;
global ( var valid, c, infile, outfile; itab);
initial infile == infile0, outfile == outfile0, c == c0;
entry c == period  $\wedge$  cvaid;
exit ((cvaid  $\wedge$ 
concat(outfile, PSI(tabrep(itab), concat(list(c), infile))) ==
concat(outfile, PHiper(tabrep(itab), list(c0), infile0)))  $\wedge$ 
outfile == concat(outfile0, PHiper(tabrep(itab), list(c0), infile0)));
V(\neg cvaid  $\wedge$ 
outfile == concat(outfile0, PHiper(tabrep(itab), list(c0), infile0)));
begin
write(outfile, Sper(c));
if eos(infile) then valid  $\leftarrow$  false else read(infile, c);
end ; %scancol %

function Sper(i: integer): token;
entry true; exit true; external ;
begin

```

```

if eof(infile) then
begin
  valid ← false;
  write(outfile, voidtoken(periodtoken));
end else
begin
  read(infile, c);
  if c = period then
begin
  write(outfile, voidtoken(periodtoken));
  if eof(infile) then valid ← false else read(infile, c)
end else write(outfile, voidtoken(periodtoken));
end
end; %scanner%
%
```

%main program%

```

entry eof(outfile) ∧
infile = infile0 ∧
¬eof(infile) ∧
itab = itab0 ∧
isitable(itab);
exit outfile = PSI(tabrep(itab0), infile0);

begin
read(infile, c);
valid ← true;
repeat
  if letter(c) then scanident else
  if digit(c) then scannumber else
  if delim(c) then scandelim else
  if c = colon then scancolon else
  if c = period then scanner else
  if eof(infile) then valid ← false else read(infile, c)
until not valid
invariant
(((valid ∧ PSI(tabrep(itab0), infile0)) =
concat(outfile, PSI(tabrep(itab0), concat(list(c), infile)))) =
∨(¬valid ∧ PSI(tabrep(itab0), infile0) = outfile)) ∧
(valid ∨ eof(infile)) ∧
isitable(itab);
end.
```

### 3. Typical verification conditions

To prove partial correctness of the scanner as presented on the previous page requires to prove 46 verification conditions. Some of these are fairly simple while others require elaborate proofs. All verification conditions can be proven by the Stanford verifier. Here are some typical examples.

Unsimplified Verification Condition: stringequal 4

$$(s_1.length \neq s_2.length \rightarrow \text{false} = (\text{stringrep}(s_1) = \text{stringrep}(s_2)))$$

Unsimplified Verification Condition: append 1

$$(s = s_0 \wedge \neg(s.length = 20) \wedge s_1 \leq s.length, s.length + 1 > \wedge s_0 = < s_1, .value, < s_1.value, [s_1.length], c >> \rightarrow \text{stringrep}(s\_0) = \text{concat}(\text{stringrep}(s_0), \text{list}(c)))$$

Unsimplified Verification Condition: scanner 2

$$(c = period \wedge valid \wedge infile = infile0 \wedge outfile = outfile0 \wedge c = 0 \wedge \neg eof(infile) \rightarrow \neg empty(infile) \wedge (infile\_2 = res(infile) \wedge c\_2 = first(infile) \wedge c\_2 = read\_x\_x(infile, c) \wedge infile\_2 = read\_f\_f(infile, c) \wedge c\_2 = period \wedge$$

```

voidoken(dperiodtoken) = mkvoidoken(dperiodtoken) ∧
outfile_2 = append(outfile, voidoken(dperiodtoken)) ∧
outfile_2 = write_f_(outfile, voidoken(dperiodtoken)) ∧
¬eo_(infile_2)
→
empty(infile_2) ∧
(infile_1 = read_f_(infile_2, c_2) ∧
c_1 = read_z_x(infile_2, c_2) ∧
infile_1 = rest(infile_2)
→
(valid ∧
concat(outfile_2, pos(tabrep(itab), concat(list(c_-1), infile_1))) =
concat(outfile_0, phiper(tabrep(itab), list(c_0), infile_0)) ∨
¬valid ∧
outfile_2 = concat(outfile_0, phiper(tabrep(itab), list(c_0), infile_0)) ∧
(valid ∨
eof(infile_1)))

```

## Unimplified Verification Condition: main 11

```

(eo_(outfile) ∧
infile = infile_0 ∧
eof(infile) ∧
itab = itab0 ∧
isitable(itab)
→
empty(infile) ∧
(infile_8 = rest(infile) ∧
c_8 = first(infile) ∧
c_8 = read_z_x(infile, c) ∧
infile_8 = read_f_(infile, c) ∧
letter(c_8) ∧
digit(c_8) ∧
delim(c_8) ∧
c_8 = colon
→
true ∧
c_8 = colon ∧
(outfile_4 = scancol_outfile(true, c_8, infile_8, outfile, itab) ∧
c_7 = scancol_c(true, c_8, infile_8, outfile, itab) ∧

```

```

valid_4 := scancol_c(valid,true,c_8,infile_8,outfile,itab) ∧
infile_7 := scancol_infile(true,c_8,infile_8,outfile,itab) ∧
(valid_4 ∧
concat(outfile_4,ps(tabrep(itab),concat(list(c_7),infile_7))) =
concat(outfile,phiper(tabrep(itab),list(c_8),infile_8)) ∨
¬valid_4 ∧
outfile_4 = concat(outfile,phiper(tabrep(itab),list(c_8),infile_8))) ∧
(valid_4 ∨
eof(infile_7))
→
isitable(itab) ∧
(valid_4 ∧
ps(tabrep(itab),infile_0) =
concat(outfile_4,ps(tabrep(itab),concat(list(c_7),infile_7))) ∨
¬valid_4 ∧
ps(tabrep(itab),infile_0) = outfile_4) ∧
(valid_4 ∨
eof(infile_7)) ∧
(isitable(itab_0) ∧
(valid_0 ∧
ps(tabrep(itab_0),infile_0) =
concat(outfile_0,ps(tabrep(itab_0),concat(list(c_0),infile_0))) ∨
¬valid_0 ∧
ps(tabrep(itab_0),infile_0) = outfile_4) ∧
(valid_0 ∨
eof(infile_0)) ∨
¬¬valid_0
→
outfile_0 = ps(tabrep(itab_0),infile_0)))

```

- Where  $\text{production}(l)$  means  $P_l$  and  $\text{rhs}$  selects the right hand side of a production  $\langle l, r \rangle$ .
- Definition of slrrel**
- $$\begin{aligned} \text{slrrel}(\text{initial\_state}, \langle \rangle) \\ \text{slrrel}(u, v) \wedge s = \text{slr}(\text{hd}(\text{last}(1, u)), n). \text{state} \\ \implies \text{slrrel}(u \cdot s, v \cdot \{n\}) \end{aligned}$$
- A derived lemma**
- $$\text{slrrel}(\text{remain}(n, s), \text{remain}(n, nt))$$

In this section we give the theory necessary to verify the parser. Several details (e.g. about sequences) are omitted.

### 1.1. Representation functions

#### 1.2. LR theory

All parsing actions are distinct

$$\begin{aligned} \text{accept} &\neq \text{reduce} \wedge \\ \text{accept} &\neq \text{shift/reduce} \wedge \\ \text{accept} &\neq \text{error} \wedge \\ \text{accept} &\neq \text{shift} \wedge \\ \text{reduce} &\neq \text{shift/reduce} \wedge \\ \text{reduce} &\neq \text{error} \wedge \\ \text{shift/reduce} &\neq \text{error} \wedge \\ \text{shift/reduce} &\neq \text{shift} \wedge \\ \text{shift} &\neq \text{error}; \end{aligned}$$

$$x \neq \text{error} \wedge x \neq \text{shift} \wedge x \neq \text{reduce} \wedge x \neq \text{shift/reduce} \implies x = \text{accept}$$

**Definition of isderiv**

$$\begin{aligned} \text{isderiv}(\text{mkforest}(x), x) \\ \text{isderiv}(u \cdot s \cdot v, w) \wedge \text{roots}(u \cdot s) = \text{rhs}(\text{production}(l)) \\ \implies \text{isderiv}(v \cdot \{l\} \cdot v, w) \\ \text{isderiv}(v, w) \wedge \{l\} = \text{rhs}(\text{production}(l)) \\ \implies \text{isderiv}(\{l\} \cdot v, w) \end{aligned}$$

Where  $\text{production}(l)$  means  $P_l$  and  $\text{rhs}$  selects the right hand side of a production  $\langle l, r \rangle$ .

**Definition of slrrel**

$$\begin{aligned} \text{slrrel}(\text{initial\_state}, \langle \rangle) \\ \text{slrrel}(u, v) \wedge s = \text{slr}(\text{hd}(\text{last}(1, u)), n). \text{state} \\ \implies \text{slrrel}(u \cdot s, v \cdot \{n\}) \end{aligned}$$

**A derived lemma**

$$\text{slrrel}(\text{remain}(n, s), \text{remain}(n, nt))$$

**Properties of LR-parsing tables**

$$\begin{aligned} \text{slrrel}(st, nt) \wedge \text{slr}(s, x).skind = \text{reduce} \wedge \text{list}(s) = \text{last}(1, st) \\ \implies \text{last}(\text{length}(\text{slr}(s, x).prod), nt) = \text{rhs}(\text{slr}(s, x).prod) \\ \text{slrrel}(st, nt) \wedge \text{slr}(s, x).skind = \text{shift/reduce} \wedge s = \text{hd}(\text{last}(1, st)) \\ \implies \text{last}(\text{length}(\text{slr}(s, x).prod) - 1, nt).list(x) = \text{rhs}(\text{slr}(s, x).prod) \\ \text{slrrel}(st, nt) \wedge \text{slr}(\text{hd}(\text{last}(1, st)), x.syn).skind = \text{accept} \\ \implies \text{nt.list}(x.syn) = \text{rhs}(\text{slr}(\text{hd}(\text{last}(1, st)), x.syn).prod) \wedge \\ \text{len}(nt) = 1 \wedge \\ \text{length}(\text{slr}(\text{hd}(\text{last}(1, st)), x.syn).prod) = 2 \wedge \\ \text{rhs}(\text{slr}(\text{hd}(\text{last}(1, st)), x.syn).prod) = \text{startsymbol} \wedge \\ x = \text{eof\_symbol} \\ \text{true} \rightarrow \text{slr}(\text{st}, \text{rhs}(p)).skind \neq \text{accept} \end{aligned}$$

### 1.3. Tree transformations

$\text{element}(x, n)$  selects element number  $n$  from sequence  $x$ . It is defined in terms of  $\text{astseqrep}$ .

$$\begin{aligned} \text{element}(\text{astseqrep}(rc, ar, s, l), n) = \text{synrep}(rc, ar[s + n - 1]) \\ \text{isderiv}(u \cdot s \cdot v, w) \wedge \text{roots}(u \cdot s) = \text{rhs}(\text{production}(l)) \\ \implies \text{isderiv}(v \cdot \{l\} \cdot v, w) \\ \text{isderiv}(v, w) \wedge \{l\} = \text{rhs}(\text{production}(l)) \\ \implies \text{isderiv}(\{l\} \cdot v, w) \end{aligned}$$

Some of the definitional clauses for  $\text{trtr}$  are given below. The remaining clauses are obvious from the definition of  $E$ .

**1.4. Extension operations**

**Definition:** subclass

```

subclass(z);
subclass(x, z ∪ p);
subclass(x, z) ∧ subclass(z, y) ==> subclass(x, y)

```

A derived lemma assuming standard interpretation of pointer-to.

```

subclass(r1, rc2) ∧ ~pointer-to(p, rcl) ==>
subclass(r1, rc1, < rc2, ⊂ p ∩, ε >)

```

**Definition of proper**

```

pointer-to(p, rc) ==> proper(synrep(rc, p))
proper(astseqrep(rc, ar, s, f)) ∧ s ≤ i ∧ i ≤ s + f - 1
      ==> proper(synrep(rc, ar[i]))
subclass(r, q) ∧ proper(synrep(r, p)) ==> proper(synrep(q, p))
empty(x) ==> proper(x)
subclass(r, q) ∧ proper(astseqrep(r, p, f, t)) ==>
proper(ar, p, f, t)
proper(synrep(ar, pt)) ∧ proper(astseqrep(rc, ar, s, f))
      ==> proper(astseqrep(rc, ar, s, f))
proper(astseqrep(rc, ar, s, f - 1)) ∧ proper(synrep(rc, ar[s + f - 1]))
      ==> proper(astseqrep(rc, ar, s, f))
proper(astseqrep(rc, ar, s, f)) ∧ subclass(rc, pc)
      ==> astseqrep(rc, ar, s, f) = astseqrep(rc, ar, s, f)
proper(astseqrep(rc, ar, s, f)) ∧ prime ≤ s ∧ f ≤ fprime
      ==> proper(astseqrep(rc, ar, s, f))

```

**2. The program**

```

trtr(prog_1, rhs) = element(rhs, 1)
trtr(prog_1, rhs) = mkeprogram(element(rhs, 2))
trtr(cdef_1, rhs) = mkenullist
trtr(cdef_1, rhs) = mkeconst(element(rhs, 2))
trtr(cdef_1, rhs) = mkelist(element(rhs, 1))
trtr(cdef_1, rhs) = mkeappend(element(rhs, 1), mkelist(element(rhs, 2)))
trtr(cde_1, rhs) = mkepair(element(rhs, 1), element(rhs, 3))
trtr(tdef_1, rhs) = mkenullist
trtr(tdef_2, rhs) = mketyp(element(rhs, 2))
trtr(tdef_1, rhs) = mkelist(element(rhs, 1))
trtr(tdef_2, rhs) = mkeappend(element(rhs, 1), mkelist(element(rhs, 2)))
trtr(tdef_1, rhs) = mkepair(element(rhs, 1), element(rhs, 3))
trtr(iden_1, rhs) = mkenum(element(rhs, 2))

type
abstract-functions = (mkedummy, mkenullist);
plabel = (termprod,
          z_1, prog_1, block_1, cdef_1,
          cdef_2, cdef_1, cdef_2, cdef_1,
          tdef_1, tdef_2, tdef_1, tdef_1,
          tdef_1, iden_1, iden_2, tden_3,
          tden_4, tden_5, tden_6, tden_7,
          idl_1, idl_2, udecl_1, udecl_2,
          udec_1, udec_2, p/deep_1,
          p/deep_2, p/deep_3, f/dec_1,
          pdec_1, ctnt_1, com_1, com_2, com_3,
          stmt_1, stmt_2, stmt_3, stmt_4,
          stmt_5, stmt_6, stmt_7, stmt_8, stmt_9,
          stmt_10, expr_1, expr_2,
          conj_1, conj_2, conj_3, rel_1, rel_2,
          rel_3, sum_1, sum_2, sum_3,
          term_1, term_2, fact_1, fact_2, fact_3,
          fact_4, fact_5, fact_6,
          vble_1, vble_2, vble_3, vble_4, parme_1,
          parme_2, parml_1, parml_2,
          parml_1, parml_2, expr_1, expr_2,
          idn_1, idn_2);

tree = 1/2;
termnonterm = (startsymbol);
token = record
  syntermnonterm;
  sem:integer;
  void:boolean end;
tokenfile = file of token;
asynkind = (tagterminal,
            tagappnd, tagarray, tagasign,
            tagunion, tagblock, tagcommandlist,
            tagconst, tagderef, tagdummy,
            tagenum, tagfcall, tagfunction,
```

```

tagto, tagif, tagindex,
taglabel, tagtypedec, tagconstdec,
tagvardec, tagcall, tagpointer,
tagprocedure, tagprogram, tagrecord,
tagrepeat, tagselect, tagtmt,
tagwhere, tagtype, tagtypede,
tagunop, tagulp, tagvard,
tagvarp, tagwhile;

atree $\leftarrow$ ↑ anode;
anode $\leftarrow$  record skind: synkind;
sub1: atree;
sub2: atree;
sub3: atree;
sub4: atree;
next: atree;
info: integer end;
stateet $\leftarrow$  1:2;
astsequence $\leftarrow$  array [1:stacklen] of atree;
treesequence $\leftarrow$  [null, sequence];
ntsequence $\leftarrow$  [b1, b2];
statesequence $\leftarrow$  array [1:stacklen] of stateet;
action $\leftarrow$  record
skind: [accept, error, shift, reduce, shiftreduce];
prodplabel;
state: stateet end;
var
list, treestack: treesequence;
slet, statestackptr: integer;
nlist, ntstack: ntsequence;
alot, astackptr: integer;
astack: astosequence;
statestack: statesequence;
input: tokenfile;
m, n: integer;
act: action;
initial_state: stateet;
eof_symbol: token;
paretree: tree;

```

```

synTree: tree;
ptr, dat: atree;
look: token;
%Building abstract syntax trees %

procedure Cmkfunction(p1, p2, p3, p4: atree; var result: atree);
global ( var #anode);
initial #anode  $\leftarrow$  anode0;
entry proper(synrep[#anode, p1])  $\wedge$  proper(synrep[#anode, p2])  $\wedge$ 
proper(synrep[#anode, p3])  $\wedge$  proper(synrep[#anode, p4]);
exit mkfunction(synrep[#anode, p1], synrep[#anode, p2],
synrep[#anode, p3], synrep[#anode, p4])  $\leftarrow$  synrep[#anode, result]  $\wedge$ 
subclass(anode0, #anode)  $\wedge$ 
proper(synrep[#anode, result]);
begin
new(result);
result↑ .skind  $\leftarrow$  tagfunction;
result↑ .subt1  $\leftarrow$  p1;
comment proper(synrep[#anode, result])  $\wedge$ 
proper(synrep[#anode, p1])  $\wedge$  proper(synrep[#anode, p2])  $\wedge$ 
proper(synrep[#anode, p3])  $\wedge$  proper(synrep[#anode, p4])  $\wedge$ 
result↑ .skind = tagfunction  $\wedge$  result↑ .subt1 = p1  $\wedge$ 
subclass(anode0, #anode);
result↑ .subt2  $\leftarrow$  p2;
result↑ .subt3  $\leftarrow$  p3;
comment proper(synrep[#anode, result])  $\wedge$ 
result↑ .skind = tagfunction  $\wedge$  result↑ .subt1 = p1  $\wedge$ 
result↑ .subt2 = p2  $\wedge$ 
result↑ .subt3 = p3  $\wedge$  proper(synrep[#anode, p4])  $\wedge$ 
subclass(anode0, #anode);
result↑ .subt4  $\leftarrow$  p4;
end;
...;
< one function for each clause of abstract syntax >
...
```

```

procedure Cmkappend(p1, p2: atree; var result: atree);
global ( var #anode);

```

```

initial #anode = anode0;
entry proper(synrep[#anode, p1]) ∧ proper(synrep[#anode, p2]) ∧
    islist(synrep[#anode, p1]) ∧ islist(synrep[#anode, p2]);
exit  mkeappend(synrep[#anode, p1], synrep[#anode, p2]);
      = synrep[#anode, result] ∧
      subclass(anode0, #anode) ∧
      islist(synrep[#anode, result]) ∧
      proper(synrep[#anode, result]);
var   t1: atree;
begin
  if p1 = nil then result ← p2 else
    begin
      new(result);
      result ↑ .subt1 ← p1 ↑ .subt1;
      comment subclass(anode0, #anode);
      Cmkappend(p1 ↑ .next, p2, t1);
      result ↑ .next ← t1;
      comment synrep[#anode, #anode ⊂ p1 ⊂ .next] =
          selrest(synrep[#anode, p1]) ∧
          isnulllist(synrep[#anode, p1]) = truthrep(false) ∧
          synrep[#anode, #anode ⊂ p1 ⊂ .subt1] =
          selfsel(synrep[#anode, p1]);
    end;
  end ;
end;

result ↑ .subt2 ← p2;
end;

%ctrfr implements ctr %
procedure ctrfr(p: plabel; tstart, tlen: integer; var result: atree);
global ( var #anode; aststack);
initial #anode = #anode0;
entry len(astsegregr[#anode, aststack, tstart, tlen]) = length(p) ∧
proper(astsegregr[#anode, aststack, tstart, tlen]);
exit  synrep[#anode, result] =
      ctrfr(p, astsegregr[#anode, aststack, tstart, tlen]) ∧
proper(synrep[#anode, result]) ∧
subclass(#anode0, #anode);
begin
  var t1, t2: atree;
  ...
  case p of
    ...
    z-1: %z ::= prog efsymbol %
      result ← aststack[tstart];
      ...
      < one case for each label of G >
      ...
  ...
stmt-5: %stmt ::= ifsymbol expr thensymbol com flaysymbol %
begin
  Cmkldummy[t1];
  Cmkaatmt[t1, t2];
  Cmkeif(aststack[tstart + 1], aststack[tstart + 3], t2, result);
end;

:dn-2: %idn ::= numbersymbol %
result ← aststack[tstart] end;
begin
  new(result);
  result ↑ .subt1 ← tag/call;
  result ↑ .subt1 ← p1;
end;

```

```
%&termir does tree transformations for terminals %

procedure Ctermir(t: tokens; var result: astree);
global ( var #anode);
initial #anode = #anode0;
entry true;
exit  synrep(#anode, result) = termir(t) ∧
proper(synrep(#anode, result)) ∧
subclsel#anode0, #anode);
begin
new(result);
result.t.ckind ← tagterminal;
result.t.info ← t.sem;
end;

procedure npop( var t: treesequence; n: integer; var list: treesequence);
initial t = t0;
entry true;
exit t = remain(n, t0) ∧ list = last(n, t0) ∧ t0 = concat(t, list);
external ;

procedure npop( var t: treesequence; n: integer; var list: treesequence);
initial t = t0;
entry true;
exit t = remain(n, t0) ∧ list = last(n, t0) ∧ t0 = concat(t, list);
external ;

procedure npush( var ts: treesequence; t: tree);
initial ts = ts0;
entry true;
exit ts = concat(ts0, list(t));
external ;

procedure npush( var ts: treesequence; t: termnonterm);
initial ts = ts0;
entry true;
exit ts = concat(ts0, list(t));
external ;

procedure tclear( var t: treesequence);
```

```
entry true; exit empty(t); external ;

procedure nclear( var t: treesequence);
entry true; exit empty(t); external ;

function slr(s: stateset; x: termnonterm): action;
entry true; exit true; external ;

procedure error;
entry true; exit false; external ;

function length(p: plabel): integer;
entry true; exit length ≥ 0; external ;

function lns(p: plabel): termnonterm;
entry true; exit true; external ;

function append(s: treesequence; t: tree): treesequence;
entry true; exit append = concat(s, list(t)); external ;

function mktree(p: plabel,
nt: termnonterm;
t: treesequence;
sm: integer): tree;
entry true; exit true; external ;

procedure apush(var s: integer; z: tree);
global ( var aststack, #anode);
initial s = s0, aststack = aststack0, #anode = #anode0;
entry proper(astseqrep(#anode, aststack, 1, s)) ∧ proper(synrep(#anode, z));
exit astseqrep(#anode, aststack, 1, s) =
concat(astseqrep(#anode, aststack0, 1, s0), list(synrep(#anode, z))) ∧
proper(astseqrep(#anode, aststack, 1, s));
begin
s ← s + 1;
aststack[s] ← z;
end;

%main parsing loop %

entry concat(input, list eof - symbol)) = input0 ∧
```

```

empty(stateerep(statestack, 1, statestackptr)) ∧
uniqueof[input] ∧
empty(asterep[#anode, aststack, 1, aststackptr]);

$$\text{kit } \text{isderiv}(\text{list}[\text{parsetree}], \text{input}[] \wedge \text{root}[\text{parsetree}]) = \text{startsymbol} \wedge$$

tree[parsetree] = synrep[#anode, act];

$$\begin{aligned} &\text{begin} \\ &\quad \text{tclear(treestack)}; \\ &\quad \text{nclear(ntstack)}; \\ &\quad \text{push(ntstack)}; \\ &\quad \text{if eos[input] then look } \leftarrow \text{eof\_symbol else read[input, look];} \\ &\quad \text{assert sllr[stateerep(statestack, 1, statestackptr), ntstack]} \wedge \\ &\quad \text{isderiv(mkefors(eqrep[look, input]), input0) \wedge} \\ &\quad \text{tree[treestack] = asterep[#anode, aststack, 1, aststackptr]} \wedge \\ &\quad \text{root[treestack]} = \text{ntstack} \wedge \\ &\quad \text{empty(ntstack)} \wedge \\ &\quad \text{stateerep(statestack, 1, statestackptr) = list(initial.state)} \wedge \\ &\quad \text{input0 = egrp[look, input]} \wedge \\ &\quad \text{unigree0[input] \wedge empty(treestack)} \wedge \\ &\quad \text{pro\_er[asterep[#anode, aststack, 1, aststackptr]]} \wedge \\ &\quad \text{len(ntstack)} = \text{len(ntstack)}; \\ &\quad \text{repeat} \\ &\quad \quad \text{act } \leftarrow \text{slr}(statestackptr, look.syn); \\ &\quad \quad \text{Ctermr(look, ptr); } \\ &\quad \quad \text{synree } \leftarrow \text{mktree(termprod, look.syn, null.sequence, look.em); } \\ &\quad \quad \text{nt } \leftarrow \text{look.syn}; \\ &\quad \quad \text{act } \leftarrow \text{slr}(statestackptr, look.syn); \\ &\quad \quad \text{Ctermr(look, ptr); } \\ &\quad \quad \text{synree } \leftarrow \text{mktree(termprod, look.syn, null.sequence, look.em); } \\ &\quad \quad \text{nt } \leftarrow \text{look.syn}; \\ &\quad \quad \text{if act.kind } \leftarrow \text{error then error else} \\ &\quad \quad \text{if act.kind } \leftarrow \text{shif then} \\ &\quad \quad \text{begin} \\ &\quad \quad \quad \text{spush(statestackptr, act.state); } \\ &\quad \quad \quad \text{npush(ntstack, look.syn); } \\ &\quad \quad \quad \text{apush(ntstack, ptr); } \\ &\quad \quad \quad \text{tpush(ntstack, synree); } \\ &\quad \quad \quad \text{get[look]; } \\ &\quad \quad \quad \text{end} \\ &\quad \quad \text{else} \\ &\quad \quad \text{if act.kind } \neq \text{accept then begin} \\ &\quad \quad \quad \text{if act.kind } \leftarrow \text{reduce then} \\ &\quad \quad \quad \text{begin} \\ &\quad \quad \quad \quad n \leftarrow \text{length(act.prod); } \\ &\quad \quad \quad \quad nt \leftarrow \text{lhs(act.prod); } \\ &\quad \quad \quad \quad \text{inpop(ntstack, n - 1, nt); } \\ &\quad \quad \quad \quad \text{anpop(aststackptr, n - 1, act); } \\ &\quad \quad \quad \quad npop(ntstack, n - 1, nlist); \\ &\quad \quad \quad \quad snpop(statestackptr, n - 1); \\ &\quad \quad \quad \text{comment} \\ &\quad \quad \quad \text{concat(nlist, list(nt)) } = \text{rhs(act.prod)} \wedge \end{aligned}$$


```

```

strell(stategrep(statestack, 1, statestackptr), nstack) &
root(nt) = nstack &
tree(nt) = asteregrp(#anode, aststack, alst, n - 1) &
len(astsegrep[#anode, aststack, alst, n - 1]) =
length(act.prod) - 1 &
len(astsegrep[#anode, aststack, alst, n - 1]) = len(nt) &
len(nt) = len(nt) &
isderiv(concat(concat(treestack, nt),
list(synree)),,
makeforest(seqrp(look, input))),,
input0) &
tree0(concat([list, list(synree)]) =
concat(astsegrep[#anode, aststack, alst, n - 1]),
list(synrep[#anode, ptr]));;
append(alst, n - 1, ptr, 0);
ctrlr(act.prod, alst, l, ptr);
nt ← lhe(act.prod);
synree ← makefref(act.prod, nt, append([list, synree], 0));
comment
tree0(synree) = synrep(#anode, ptr);
act ← slr(stop(statestackptr), nt);
end ;
tpush(treestack, synree);
apush(aststackptr, ptr);
npush(nstack, nt);
spush(statestackptr, act.state);
end ;
until act.statekind = accept
invariant
uniquef(input) &
isderiv(concat(treestack, makeforest(seqrp(look, input))), input) &
strell(stategrep(statestack, 1, statestackptr), nstack) &
tree(nt) = asteregrp(#anode, aststack, 1, aststackptr) &
proper(asteregrp(#anode, aststack, 1, aststackptr)) &
root(nt) = nstack &
len(nt) = len(nt) &
(act.statekind = accept →
(act = slr(hd([alst]), stategrep(statestack, 1, statestackptr))), nt) &
nt = look.syn);
begin
n ← length(act.prod);
nt ← lhe(act.prod);

```

```

tmpop(treestack, n - 1, nt);
anpop(aststackptr, n - 1, alst);
nnpop(nstack, n - 1, nt);
anpop(statestackptr, n - 1);
assert concat(nt, list(look.syn)) = rhs(act.prod) &
segrp(look, input) = list(leaf.symbol) &
strell(stategrep(statestack, 1, statestackptr), nstack) &
len(treestack) = len(nt) &
roots(nt) = nt & len(treestack) = 0 &
len(astsegrep(#anode, aststack, alst, n - 1)) = length(act.prod) - 1 &
len(astsegrep(#anode, aststack, 1, aststackptr)) = 0 &
trees(nt) = astsegrep(#anode, aststack, 1, aststackptr) &
nt = startsymbol &
proper(astsegrep(#anode, aststack, alst, n - 1)) &
isderiv(concat(nt,
makefrest([list, look])),,
input0);
Ctermtr(look, ptr);
append(alst, n - 1, ptr, 0);
ctrlr(act.prod, alst, l, ast);
parsetree ← makefref(act.prod, nt,
append(nt,
mktree(termprod, look.syn, null.sequence, look.syn)), 0);
end .

```

$\text{apply}[(\zeta, i)] = \text{mkfunctype}(\mu\text{list}, \tau\text{au}) \wedge$
$\tau\text{T} = \text{passablestar}(\mu\text{list}, \text{Astar}(\text{Elist}, \zeta)) \text{ by } \text{mkvalmode}(\tau\text{au});$
<b>sem40:</b> replace $\text{Ae}[\text{mkfuncall}(I, \text{Elist}), \zeta]$ where
$\tau\text{T} = \text{ispfunctype}(\text{apply}[(\zeta, I)]) \wedge$
$\text{apply}[(\zeta, i)] = \text{mkfunctype}(\mu\text{list}, \tau\text{au}) \wedge$
$\tau\text{T} = \text{passablestar}(\mu\text{list}, \text{Astar}(\text{Elist}, \zeta)) \text{ by } \text{mkvalmode}(\tau\text{au});$
<b>sem41:</b> replace $\text{Ae}[\text{mkfuncall}(I, \text{Elist}), \zeta]$ where
$\tau\text{T} = \text{isfunctype}(\text{apply}[(\zeta, I)]) \text{ by } \text{Ao}[I, \text{Astar}(\text{Elist}, \zeta)];$
<b>sem42:</b> replace $\text{Ac}[\text{mkselect}(E, I), \zeta]$ where
$\text{Ac}[E, \zeta] = \text{mkvarmode}(\text{mkrecordtype}(\mu\text{u}, P\text{st})) \wedge$
$\tau\text{T} = \text{islement}(\text{mkpair}(I, \mu\text{u}), P\text{st}) \text{ by } \mu\text{u};$
<b>sem43:</b> replace $\text{Ac}[\text{mkselect}(E, I), \zeta]$ where
$\text{Ac}[E, \zeta] = \text{isvarmode}(\text{Ac}[E, \zeta]) \wedge$
$\text{Ac}[E, \zeta] = \text{mkvarmode}(\text{mkrecordtype}(\mu\text{u}, P\text{st})) \wedge$
$\tau\text{T} = \text{islement}(\text{mkpair}(I, \mu\text{u}), P\text{st}) \text{ by } \mu\text{u};$
<b>sem44:</b> replace $\text{Ac}[\text{mkindex}(E0, E1), \zeta]$ where
$\text{Atyp}[E0, \zeta] = \text{mearraytype}(\mu\text{u}, \tau\text{au1}, \tau\text{au2}) \wedge$
$\tau\text{T} = \text{compatible}(\tau\text{au1}, \text{Atyp}[E1, \zeta]) \wedge$
$\tau\text{T} = \text{isvar}(\text{Ac}[E0, \zeta]) \text{ by } \text{mkvarmode}(\tau\text{au2});$
<b>sem45:</b> replace $\text{Ac}[\text{mkindex}(E0, E1), \zeta]$ where
$\text{Atyp}[E0, \zeta] = \text{mearraytype}(\mu\text{u}, \tau\text{au1}, \tau\text{au2}) \wedge$
$\tau\text{T} = \text{compatible}(\tau\text{au1}, (\text{Atyp}[E1, \zeta])) \wedge$
$\text{FF} = \text{isvar}(\text{Ac}[E0, \zeta]) \wedge$
$\tau\text{T} = \text{isvar}(\text{Ac}[E0, \zeta]) \text{ by } \text{mkvalmode}(\tau\text{au2});$
<b>sem46:</b> replace $\text{Ac}[\text{mkderef}(E), \zeta]$ where
$\text{Atyp}[E, \zeta] = \text{mkpointertype}(\mu\text{u}, \tau\text{au})$
by $\text{mkvarmode}(\tau\text{au});$
<b>sem47:</b> replace $\text{Astar}(\text{Elist}, \zeta)$ by
$\text{mkappend}(\text{mklist}(\text{Aelset}[\text{irat}(\text{Elist})], \zeta));$
$\text{Astar}(\text{selest}(\text{Elist}), \zeta);$
<b>sem48:</b> replace $\text{Astar}(\zeta, \zeta)$
where $\text{isnullist}(\zeta) = \text{TT}$ by $\text{mknulllist};$
<b>sem37:</b> replace $\text{Ae}[\text{mkunop}(O, E), \zeta]$ by $\text{Ao}(O, \text{Ac}[E, \zeta]);$
<b>sem38:</b> replace $\text{Ae}[\text{mkbinop}(B0, \Omega\text{mega}, E1), \zeta]$ by
$\text{Aw}(\Omega\text{mega}, \text{Ac}[E0, \zeta], \text{Ac}[E1, \zeta]);$
<b>sem39:</b> replace $\text{Ae}[\text{mkfuncall}(I, \text{Elist}), \zeta]$ where
$\tau\text{T} = \text{isfunctype}(\text{apply}[(\zeta, I)]) \wedge$

### Appendix 5. Static semantics

#### 1. Logical basis

In this section we give some of the rules necessary for the verification of the static semantics. We first present the formal definition of  $\epsilon$  in the verifier's rule language. The translation of other definitions into machine readable form follows analogously.

In addition to the definition of all semantics functions operationalization of the fixed point defining recursive types is required for the proof of the static semantics. This theory is also presented below.

#### 1.1. Rules for $\epsilon$

<b>rulefile (static-<math>\epsilon</math>)</b>
constant $\tau\text{T}, \text{FF}, \text{UU}, \text{int}, \text{bool};$
<b>sem32:</b> replace $\text{Ae}[\text{mknumber}(n), \zeta]$ by
$\text{mkconstmode}(\text{An}(n), \text{mkusertype}(\text{tagrep}(\text{int})), \text{An}(n), \text{An}(n));$
<b>sem33:</b> replace $\text{Ae}[\text{mkzeroid}(I), \zeta]$ where
$\text{FF} = \text{isgpmode}(\text{apply}[(\zeta, I)]) \wedge$
$\text{FF} = \text{isprocmode}(\text{apply}[(\zeta, I)]) \wedge$
$\text{FF} = \text{isprocmode}(\text{apply}[(\zeta, I)]) \text{ by } \text{applyl}[(\zeta, I)];$
<b>sem37:</b> replace $\text{Ae}[\text{mkunop}(O, E), \zeta]$ by $\text{Ao}(O, \text{Ac}[E, \zeta]);$
<b>sem38:</b> replace $\text{Ae}[\text{mkbinop}(B0, \Omega\text{mega}, E1), \zeta]$ by
$\text{Aw}(\Omega\text{mega}, \text{Ac}[E0, \zeta], \text{Ac}[E1, \zeta]);$
<b>sem39:</b> replace $\text{Ae}[\text{mkfuncall}(I, \text{Elist}), \zeta]$ where
$\tau\text{T} = \text{isfunctype}(\text{apply}[(\zeta, I)]) \wedge$

### 1.2. Recursive types

The function symbol  $\text{At}$  denotes the semantics function  $\mathfrak{f}$ ;  $Ft$  denotes the function  $t_f$  used in the operationalization. We first present the definition of these two functions in the rule language and then present lemmas dealing with least fixed points.

#### 1.2.1. Types

**rulefile (static-t)**

% static semantics for type %

```
constant TT, FF, UU;

st1: replace At(mktypeid(I), zeta0, zeta1) where
      apply1(zeta0, I) = mktypeode(tau) by
      mkepair(tau, zeta0);
```

```
F st1: replace Ft(mktypeid(I), zeta0, ur) where
      apply1(zeta0, I) = mktypeode(tau) by
      mkepair(tau, zeta0, ur);
```

```
F st2: replace At(mkneuronum(Ist), zeta0, zeta1) where
      mkepair(zeta3, n) = etstar(Ist, nu, zeta2, intrep(0)) ∧
      newtag(zeta0) = mkepair(nu, zeta2) by
      mkepair(mknesubtype(nu, intrep(1), n), zeta3);
```

%Enumeration %

```
st2: replace At(mkneuronum(Ist), zeta0, zeta1) where
      mkepair(zeta3, n) = etstar(Ist, nu, zeta2, intrep(0)) ∧
      newtag(zeta0) = mkepair(nu, zeta2) by
      mkepair(mknesubtype(nu, intrep(1), n), zeta3);
```

```
F st2: replace Ft(mkneuronum(Ist), zeta0, ur) where
      mkepair(zeta3, n) = etstar(Ist, nu, zeta2, intrep(0)) ∧
      newtag(zeta0) = mkepair(nu, zeta2) by
      mkepair(mknesubtype(nu, intrep(1), n), zeta3, ur);
```

%where we introduced and auxiliary function etstar as defined below: %

```
et1: replace etstar(Ist, nu, zeta, n) where isnull(Ist) = TT
      by mkepair(zeta, n);

et2: replace etstar(Ist, nu, zeta, n) where isnull(Ist) = FF
      by etstar(selrest(Ist), nu, et(selfrst(Ist), nu, zeta, succ(n))),
```

```
et3: replace et(id, nu, zeta, n) by
      redefine(zeta, mkconsatmode(n, mkectype(nu, n, n)), id);
succ1: replace succ(intrep(n)) by intrep(n + 1);

st3: replace At(mkessubrange(E1, E2), zeta0, zeta1) where
      isconstmode(Ae(E1, zeta0)) = TT ∧
      isconstmode(Ae(E2, zeta0)) = TT by
      mkepair(Aunion(Atype(Ae(E1, zeta0)), Atype(Ae(E2, zeta0))), zeta0);

F st3: replace Ft(mkessubrange(E1, E2), zeta0, ur) where
      isconstmode(Ae(E1, zeta0)) = TT ∧
      isconstmode(Ae(E2, zeta0)) = TT by
      mkepair(Aunion(Atype(Ae(E1, zeta0)), Atype(Ae(E2, zeta0)))), zeta0, ur);

st4: replace At(mkearray(T1, T2), zeta0, zeta1) where
      At(T1, zeta0, zeta1) = mkepair(tau1, zeta2) ∧
      At(T2, zeta2, zeta1) = mkepair(tau2, zeta3) ∧
      isstype(tau1) = TT ∧
      newtag(zeta3) = mkepair(nu, zeta4) by
      mkepair(mkearraytype(nu, tau1, tau2), zeta4);

F st4: replace Ft(mkearray(T1, T2), zeta0, ur) where
      F(T1, zeta0, ur) = mketriplet(tau1, zeta2, ur2) ∧
      F(T2, zeta2, ur2) = mketriplet(tau2, zeta3, ur3) ∧
      isstype(tau1) = TT ∧
      newtag(zeta3) = mkepair(nu, zeta4) by
      mketriplet(mkearraytype(nu, tau1, tau2), zeta4, ur3);

st5: replace At(mkerecord(IList), zeta0, zeta1) where
      distinct(selidist(IList)) = TT ∧
      mkepair(tau1, zeta2) = Atstar(IList, zeta0, zeta1) ∧
      mkepair(nu, zeta3) - newtag(zeta2) by
      mkepair(mkerecordtype(nu, mkeccoc(selidist(IList), tau1))), zeta3;
```

```
F st5: replace Ft(mkerecord(IList), zeta0, ur) where
      distinct(selidist(IList)) = TT ∧
      succ(n));
```

```

mke triple(zeta1, zeta2, ur1) = Fstar([Tlist, zeta0, ur]  $\wedge$ 
  mke pair(nu, zeta3) = newtag(zeta2) by
  mke triple(
    mke recordtype(nu, mke record(zeta1, tau1)),
    zeta3, ur1);

```

```

st6: replace At(mke pointer(l), zeta0, zeta1) where
  newtag(zeta0) = mke pair(nu, atype(apply(zeta1, l))), zeta2);

F46: replace F1(mke pointer(l), zeta, ur) where
  newtag(zeta) = mke pair(nu, zeta1)  $\wedge$ 
  mke triple(mke pointer type(nu, tau), zeta2, ur2);

star1: replace Atstar(dl, z1, z2) where isnull(dl) = TT by z1;
star2: replace Atstar(dl, z1, z2) where
  mke pair(tau, z3) = At(isel(isrl(dl), z1, z2)) BY
  mke pair(tau, z3) = At(isel(isrl(dl), z1, z2));

```

### 1.2.2. Fixed points

#### rulefile (fixedpoints)

constant emptyset, emptyur;

%definition of resolve %

```

fix1: replace resolve(a, m, t, e, z, ur)
  where -isnull(ur)  $\wedge$ 
  apply(enurep(a, m, t, e, z, ur)) =
  mke type mode((typerep(t, tau)))
  by resolve(a, m, <, t, C tpos(ur) D, t C tau D>, e, z, neztos(ur));

```

```

fix2: replace resolve(a, m, t, e, z, ur) where isnull(ur) by t;

```

%resolve computes a least fixed point %

```

fix3: replace /iz(enurep(a, m, t, e, z, ur)) by
  enurep(a, m, resolve(a, m, t, e, z, ur)), e, z);

```

### 2. The program

The development of the program from specifications is straightforward. The refinement step for *Ce* is described in the text in some detail. Below we give a listing of the declarations for the static semantics implementation, the code of *Ce*, and routines used to compute recursive types.

#### 2.1. Declarations

##### 2.1.1. Types

pascal

```

type termnonterm = (tagindex, tagselect, taglabel, tagstmt,
  tagcommandlist, tagtype, tagconst, tagvard,
  tagconstdec, tagtypedec, tagvardec, tagfunction,
  tagprocedure, tagarray, tagvar, tagqual, tagenum,
  tagtyped, tagarray, tagrecord, tagrange,
  tagpointer, tagprogram, tagif, tagassign,
  tagwhile, tagrepeat, taggoto, tagcall,
  tagblock, tagdummy, tagespid, tagnumber,
  tagunop, tagbinop, tagderec, tagcast);

atree =  $\uparrow$  anode;
anode = record
  skind: termnonterm;
  sub1: atree;
  sub2: atree;
  sub3: atree;
  sub4: atree;
  next: atree;
  info: integer
end;

ttype =  $\uparrow$  tnode;
mode =  $\uparrow$  mnode;
tnode = record
  tknd: (tagrecordtype, tagarraytype,
  tagsubtype, tagpointertype,
  tagnitype);
  typetag: integer;

```

```

lwb: integer;
upb: integer;
sub1: ttype;
sub2: ttype;
recs: type;
id: integer
end;

mnode = record
  mkind: (tagvalmode, tagupemode, tagfuncmode,
  tagsfuncmode, tagarmode, tagypemode,
  tagprocmode, tagsprocmode, tagfunkmode,
  tagconstmode);
  ty: ttype;
  misit: mode;
  next: mode;
  val: integer
end;

static_environment = ↑ enode;
enode = record
  id: atree;
  md: mode;
  next: static_environment
end;

nullary_functions = (TT, FF, int, bool);

uptr = ↑ unode;
unode = record tp: ttype; id: atree; next: uptr end;

2.1.2. Abstract syntax

%Testing abstract syntax %

function Cisindex(e:atree):boolean;
global (#anode);
entry true;
exit truthrep(Cisindex) = isindex(symrep(#anode,e));
begin
  Cisindex ← e ↑ .skind = tagindex;
end;

```

```

function Ciselect(e:atree):boolean;
global (#anode);
entry true;
exit truthrep(Ciselect) = isselect(symrep(#anode,e));
begin
  Ciselect ← e ↑ .skind = tagselect;
end;

%Decomposition of abstract syntax %

procedure matchderej(e:atree; var el:atree);
global (#anode);
entry isderej(symrep(#anode,e)) = TT;
exit makederej(symrep(#anode,e)) = symrep(#anode,e);
begin
  el ← e ↑ .sub1;
end;

procedure matchcall(e:atree; var i:atree);
global (#anode);
entry isfcall(symrep(#anode,e)) = TT;
exit makedefcall(symrep(#anode,i),symrep(#anode,e)) = symrep(#anode,e);
begin
  i ← e ↑ .sub1;
  elist ← e ↑ .sub2;
end;

%Testing types %

function Cisarraytype(tau: ttype): boolean;
global (#tnode);
entry true;
exit isarraytype(symrep(#tnode,tau)) = truthtable(Cisarraytype);
begin
  Cisarraytype ← tau ↑ .tkind = tagarraytype;
end;

```

```

Cisafuncnode  $\leftarrow$  mu  $\uparrow$  .mkind = tagofunnode;
end;

%constructing new types %

procedure Cmkearraytype(nu: integer; tau1, tau2: ttype; var r: ttype);
global (#mnode, #tnode);
exit typerp(#tnode, r) =
    mkearraytype(typerp(nu), typerp(#tnode, tau1), typerp(#tnode, tau2));
begin
new(r);
r  $\uparrow$  .typetag  $\leftarrow$  nu;
r  $\uparrow$  .tkind  $\leftarrow$  tagarraytype;
r  $\uparrow$  .subt1  $\leftarrow$  tau1;
r  $\uparrow$  .subt2  $\leftarrow$  tau2;
end;
...

%decomposing types %

procedure matcharraytype(tau: ttype; var nu: integer; var tau1, tau2: ttype);
global (#tnode);
entry isarraytype(typerp(#tnode, tau)) = TT;
exit typerp(#tnode, tau) =
    mkearraytype(typerp(nu), typerp(#tnode, tau1), typerp(#tnode, tau2));
begin
nu  $\leftarrow$  tau  $\uparrow$  .typetag;
tau1  $\leftarrow$  tau  $\uparrow$  .subt1;
tau2  $\leftarrow$  tau  $\uparrow$  .subt2;
end;
...

%testing modes %

function Cisafuncnode(mu: mode): boolean;
global (#tnode, #mnode);
entry true;
exit isafuncnode(maderep(#mnode, #tnode, mu)) =
    truthrep(Cisafuncnode);
begin

```

%constructing modes %

```

procedure Cmkfuncnode(mu: mode; tau: ttype; var mu: mode);
global (#tnode, #mnode, #anode);
entry true;
exit maderep(#mnode, #tnode, mu) =
    mkefuncnode(maderep(#mnode, #mnode, #tnode, mu), typerp(#tnode, tau));
begin
new(mu);
mu  $\uparrow$  .mkind  $\leftarrow$  tagofunnode;
mu  $\uparrow$  .mlist  $\leftarrow$  mu;
mu  $\uparrow$  .ty  $\leftarrow$  tau;
end;
...

%decomposing modes %

procedure matchfuncnode(mu: mode; var mlist: mode; var tau: ttype);
global (#mnode, #tnode);
entry isafuncnode(maderep(#mnode, #tnode, mu)) = TT;
exit maderep(#mnode, #tnode, mu) =
    mkefuncnode(maderep(#mnode, #mnode, #tnode, mu)) =
        typerp(#tnode, tau));
begin
tau  $\leftarrow$  mu  $\uparrow$  .ty;
mlist  $\leftarrow$  mu  $\uparrow$  .mlist;
end;
...

2.1.3. Auxiliary functions
procedure error;
entry true; exit false; external;
function Colop(atree: mnode): mode;
global (#tnode, #mnode, #anode);

```

```

entry true;
exit moderep(#mnode, #tnode, Co) ==
  Aofsynrep(#anode, op), moderep(#mnode, #tnode, md);
external ;
...
%routines for binary operators, special functions, :special procedures%
...
function Cisint(Ty:ttype);boolean;
global (#enode, #mnode, #tnode, #anode);
entry true;
exit truthep(Cisint) == init(typerep(#tnode, Ty));
begin
  Cisint  $\leftarrow$  (Ctype tag(ty) == int);
end ;
function Cisbool(Ty:ttype);boolean;
...

```

```

entry true;
exit moderep(#mnode, #tnode, Co) ==
  Aofsynrep(#anode, op), moderep(#mnode, #tnode, md);
external ;
...
procedure Cunion(Ty1, Ty2:ttype; var r:ttype);
global (#enode, #mnode, #tnode, #anode);
exit typerep(#tnode, r) = Aunion(typerep(#tnode, Ty1), typerep(#tnode, Ty2));
var
  nu1, nu2, iota1, iota2, iota3, iota4:integer;
begin
  if Cissubtype(Ty1) and Cissubtype(Ty2) then
    begin
      match subtype(Ty1, nu1, iota1, iota2);
      match subtype(Ty2, nu2, iota3, iota4);
      if nu1 == nu2 then Cmksubtype(nu1, iota1, iota2, iota3, iota4, r)
        else error;
      end else error;
    end ;
...
procedure Cpassable(Md1, Md2:mode; var r:boolean);
global (#enode, #mnode, #tnode, #anode);
entry true;
exit truthep(r)
  == passable(moderep(#mnode, #tnode, Md1),
  moderep(#mnode, #tnode, Md2));
var
  tau1, tau2:ttype;
  t:boolean;
begin
  Ctype(Md1, tau1); Ctype(Md2, tau2);
  if Cisvar(Md1) then
    begin
      r  $\leftarrow$  Cisvar(Md2) and
      Ceqal([tau1, tau2]);
    end else
      if Cisval(Md1) then
        procedure Ccontains(Ty1, Ty2:ttype; var r:boolean);

```

```

begin
  Ccompatible(tau1, tau2, t);
  r ← Cisval(Md1) and t;
  end else begin error; r ← false; end ;
end;



### 2.2. Expressions


procedure Cestarr(exp:atree; zeta:static_environment; var m:mode);
global #enode, #mnode, #tnode, #anode;
entry true;
exit moderp(#mnode, #tnode, m) =
  Aestarr(signrep(#anode, exp),
  envrep(#anode, #mnode, #tnode, #enode, zeta));
forward ;

procedure Celeafp(atree; zeta:static_environment; var m:mode);
global (#enode, #mnode, #tnode, #anode);
entry true;
exit moderp(#mnode, #tnode, m) =
  Aesrep(#anode, exp),
  envrep(#anode, #mnode, #tnode, #enode, zeta));
var
  b : boolean;
  I : atree;
  n : atree;
  nu : integer;
  iota : integer;
  mid, m1, m2 : mode;
  op : atree;
  e1, e2 : atree;
  elist : atree;
  mult12, mult1 : mode;
  psi, tauprime, tau, tau1, tau2 : type;
begin
  if Cisnumber(exp) then
    begin
      matchnumber(exp, n);
      iota ← value(n);
      Crmksubtypelist(iota, iota, iota, tau);
      Crmkconsmodel(iota, tau, m);
    end else

```

```

 $m \leftarrow C.e[i, mule2];$ 
end else error;
end else
if  $C.iselect(exp)$  then
begin
matchselect(exp, e, i);
C(e, zeta, md);
if  $C.iavarode(md)$  then
begin
matcharrayode(md, tau);
if  $C.iarecordtype(tau)$  then
begin
matchrecordtype(tau, nu, ps);
m  $\leftarrow$  assoc(ps, i);
end else error;
end else
if  $C.iavarode(md)$  then
begin
matcharrayode(md, tau);
if  $C.iarecordtype(tau)$  then
begin
matchrecordtype(tau, nu, ps);
m  $\leftarrow$  assoc(ps, i);
end else error;
end else
if  $C.iisindex(exp)$  then
begin
matchindex(exp, e1, e2);
C(e1, zeta, m1);
C(e2, zeta, m2);
Ctype(m1, tau);
if  $C.iarraytype(tau)$  then
begin
matcharraytype(tau, nu, tau1, tau2);
Ctype(m2, tauprime);
Ccompatible(tau1, tauprime, b);
if b then
if  $C.evar(m1)$  then  $C.mkevarnode(tau2, m)$  else
if  $C.iaval(m1)$  then  $C.mkevalnode(tau2, m)$  else
error;
else error;
end else error;
end else
if  $C.iisindex(exp)$  then
begin
matchdere(exp, e);
C(e, zeta, md);
Ctype(Md, tau);
if  $C.ipointertype(tau)$  then
begin
matchpointer(type(tau, nu, tauprime));
end else error;
end else error;
end;

```

procedure  $C.estar;$   
%has forward declaration %  
var mu1, mu2:mode;

```

begin
if  $C.isnulllist(exp)$  then m  $\leftarrow$  nil else
begin
Ce(Cscifstr(exp), zeta, mu1);
Cstar(Ctelec(exp), zeta, mu2);
Cmodeappend(mu1, mu2, m);
end ;

```

### 2.3. Types

Some additional documentation not  
the specifications of the following routine.  
prove that changes to reference classes occur.

```

procedure  $CFit(ppc:atrec; zeta:static_environ;$   

var t:type; var z:static_environ);
global ( var #tnode; #tnode; #node; #ure0, #un  

initial #tnode := tnode0; ure0 := ure0; #un  

exit raketriplet(typep(#tnode, t),  

enrpp(#node, #node, #node, ure0));  

urep(#node, #node, ure0)) =  

urep(#synrp(#node, #node, tpe));

```

```

urep(unode0, #anode, ure0)) &
subclae((inode0, #tnode) &
~ptot(leftmost((urep((#unode, #anode, ure0)),
pointer((urep((#unode, #anode, ure0))), (tnode0)))))

var
n, nu: integer;
e1, e2, id, idist, idist0: atree;
mu, m1, m2: mode;
zeta3, zeta4, zeta5: static_environment;
r: atree;
rcs, tauList: type;
tau, t1, t2: type;
begin
if Cistypeid(tpe) then
begin
matchtypeid(tpe, id);
mu ← Capply1(zeta, id);
if Cistypemode(mu) then
begin
matchtypemode(mu, tau);
t ← tau;
z ← zeta;
end else error;
end else
if Cisenum(tpe) then
begin
matchenum(tpe, idist);
Cnewtag(zeta, mu, zeta3);
n ← 0;
idist0 ← idist;
zeta4 ← zeta3;
invariant etstar(synrep((#anode, idist0), tagrep(mu),
enrep((#anode, #tnode, #tnode, #tnode, zeta3), intrep(0))) =
etstar(synrep((#anode, idist), tagrep(mu),
enrep((#anode, #tnode, #tnode, #tnode, zeta4), intrep(n)))
while not Cisnullset(idist) do
begin
n ← n + 1;
Cmksubtypes((nu, n, n, tau));
Cmkconstmode((nu, tau, mu));
Cdefine(zeta, mu, Cselfrst(idist), zeta4);
idist ← Cselfrst(idist);
end else error
end else
begin
Cfstars(rs, zeta, tauList, zeta3, ure);
matchrecord(tpe, rs);
idist ← Cselfrst(rs);
if Cdistinct(idist) then
begin
Cfstars(rs, zeta, tauList, zeta3, ure);
Cnewtag(zeta3, mu, zeta5);
z ← zeta5;
end else error;
end else
if Cisrecord(tpe) then
begin
matchrecord(tpe, rs);
idist ← Cselfrst(rs);
if Cdistinct(idist) then
begin
Cfstars(rs, zeta, tauList, zeta3, ure);
Cnewtag(zeta3, mu, zeta5);
Cmkerec((idist, tauList, rcs));
Cmkrecordtype((nu, rcs, t));
end else error
end else
begin
Cfstars(rs, zeta, tauList, zeta3, ure);
Cnewtag(zeta3, mu, zeta5);
Cmkerec((idist, tauList, rcs));
Cmkrecordtype((nu, rcs, t));
end else error
end

```

```

if Cispointer(ep) then
begin
matchpointer(ep,id);
Cnewtag(zeta,nu,s);
addnewuid(tau,ure,z);
Cmkpointer(type[nu,tau,t]);
end else error;
end;

procedure CF(did:tree; zeta:static-environment;
var ure:ep; uptr);
var urep:uptr;
global ( var #tnode, #unode, #anode, #mnode );
initial #tnode = tnode0, ure = ure0, #unode = unode0;
exit mkpair(urep(#anode, #tnode, #tnode, #tnode, s),
urep(#unode, #anode, #anode, ure));
Fdit(synepi#anode, dec);
urep(#anode, #mnode, tnode) #anode, zeta),
urep(unode0, #anode, ure0)) /\
subcase(tnode0, #tnode) /\
prior(reminus(pointer(urep(#unode, #anode, ure0))), tnode0);
var i, t: tree;
tau:type;
mu:mode;
begin
if Cistypedecl(dec) then
begin
matchtypedec(dec,i,t);
CF(i, zeta, tau, z, ure);
Cmktypemode(tau, mu);
Cedefine(z, mu, i, z);
end else error;
end;

procedure Createol(z:static-environment; ure:uptr);
global ( var #tnode, #mnode, #anode, #unode );
initial #tnode = tnode0,
exit #tnode = resolve(#anode, #mnode, tnode0, #anode, z,
urep(#unode, #anode, ure))) /\

```

## Appendix 6. Code generation

**1. Logical basis**

---

```

apply: infer varthetaG(Sapply(gamma, epsilon), s, Tapply(gammaT, epsilon))
      from varthetaG(gamma, s, gammaT);

unop: infer varthetaG(unop(gamma, O), s, Tunop(gammaT, O))
      from varthetaG(gamma, s, gammaT);

binop: infer varthetaG(binop(gamma, O), s, Tbinop(gammaT, O))
      from varthetaG(gamma, s, gammaT);

addr: infer varthetaG(Addr(gamma, alpha, n), s,
      Tapply(Iaddr(gammaT, m), alphatT))
      from varthetaG(gamma, s, gammaT) ∧ n = m ∧
      smember(n, alpha, alphatT, s);

uthU1: infer smember(apply2(rho, I), apply1(rho, I),
      Tapply(y, CppU2(rho, I), I), PhiS(zeta, rho, y))
      from varthetaG(rho, zeta, y, rhoT);

rulefile (target-semantics)

constant nocode, nullist;

Macro: replace Macro(nocode, rhoT, gammaT) by gammaT;
Lcrl: replace LcrlT(nocode, rhoT, gammaT) by nullist;
Listl: replace cons(z, nullist) by mklist(z);
rulefile (commands)

constant nullist;

com1: replace Com(mklabel(N, Theta), zeta, rho, gamma)
      by Bcr(Theta, zeta, rho, gamma);

com2: replace Com(mkcommandt(Theta), zeta, rho, gamma)
      by Bcr(Theta, zeta, rho, gamma);

com3: replace Com(mkcommandlist(G0, G1), zeta, rho, gamma)
      by Ccr(G0, zeta, rho, Ccr(G1, zeta, rho, gamma));

Jscr1: replace Jscr(mkcommandlist(G0, G1), zeta, rho, gamma)
      from varthetaG(z, s, s);

```

by append(Jscr(G0, zeta, rho, Cscr(G), zeta, rho, gamma));  
 Jscr(G, zeta, rho, gamma));

Jscr2: replace Jscr(mkeestmt(S), zeta, rho, gamma) by nullist;

Jscr3: replace Jscr(mkelabel(N, S), zeta, rho, gamma)  
 by mkeist(Bscr(S, zeta, rho, gamma));

rulefile (statements)

constant mkendumy, mkesprcmode, mkenullist;

Bd01: replace Bscr(mkeassig(E0, E1), zeta, rho, gamma)  
 where ce(E0, zeta) = mkevarmode(tau)  
 by Lscr(E0, zeta, rho, Ascr(E1, zeta, mkevalmdef(tau), rho,  
 update(gamma)));

Bd02: replace Bscr(mkeif(B, G0, G1), zeta, rho, gamma)  
 where mkepair(zeta0, N0) = newlabel(zeta)  $\wedge$   
 mkepair(zeta1, N1) = newlabel(zeta0)  
 by Cscr(mkecommandlist(G0,  
 mkecommandlist(mkeestmt(mkelabel(N1)),  
 mkecommandlist(mkelabel(N0,  
 mkeestmt(mkelock(mkenullist, G1)))),  
 mkeestmt(mkedumy))))),  
 zeta1, rho, gamma);

Bd03: replace Bscr(S, zeta, rho, gamma)  
 where idummy(S)  
 by gamma;

Bd04: replace Bscr(mkewhile(E, G), zeta, rho, gamma)  
 where mkepair(zeta0, N0) = newlabel(zeta)  $\wedge$   
 mkepair(zeta1, N1) = newlabel(zeta0)  
 by Cscr(mkecommandlist(mkelabel(N0, mkenullist(E, N1)),  
 mkecommandlist(G,  
 mkecommandlist(mkeestmt(mkelabel(N0)),  
 mkelabel(N, mkedumy)))),  
 zeta1, rho, gamma);

Bd05: replace Bscr(mker spec(G, E), zeta, rho, gamma)

where mkepair(zeta0, N0) = newlabel(zeta)  
 by Cscr(mkecommandlist(mkelabel(N0, mkeblock(mkenullist, G)),  
 mkeint(mkenullist(E, N0))),  
 zeta0, rho, gamma);

Bd06: replace Bscr(mkegoto(N), zeta, rho, gamma)  
 by adjust(apply(rho, N) gamma);

Bd07: replace Bscr(mkepair(I, Elist), zeta, rho, gamma)  
 where isprcmode(apply(zeta, I))  
 by Escr(Elist, zeta, rho, Sscr(I, gamma));

Bd08: replace Bscr(mkepcall(I, Elist), zeta, rho, gamma)  
 where ~isprcmode(apply(zeta, I))  $\wedge$  mkepair(mulist, n) =  
 estar(Elist, zeta)  
 by Ascr(Elist, zeta, multist, zeta, rho, Spcall(rho, I, gamma));

Bd10: replace Bscr(mkeunlks(E0, N0), zeta, rho, gamma)  
 by Rscr(E0, zeta, rho, Cnd(gamma, apply(rho, N0)));

rulefile (expressions)

aux1: from is/unmode(x) infer x = mkefunmode;  
 constant nullist, mkefunmode;

ed01: replace Escr(E, zeta, rho, gamma)  
 where ce(E, zeta) = mkeconstmode(epsilon, tau)  
 by apply(gamma, epsilon);

ed02: replace Escr(E, zeta, rho, gamma)  
 where ~isconstmode(cc(E, zeta))  
 by AEscr(E, zeta, rho, gamma);

ed03: replace AEscr(E, zeta, rho, gamma)  
 where E = mkenumber(N) by apply(gamma, N);

ed04: replace AEscr(E, zeta, rho, gamma)  
 where E = mkeprod(I)  $\wedge$   
 apply(zeta, I) = mkefunmode  
 by Sscr(I, gamma);

ed05: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkeepid(I) \wedge$   
 isafunmode(apply(zeta, I))  
 by  $fcall[apply3(rho, I), gamma]$ ;

ed06: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkeepid(I) \wedge$   
 ispfunmode(apply(zeta, I))  
 by  $fcall[apply3(rho, I), gamma]$ ;

ed07: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkeepid(I) \wedge$   
 "is; unmode(apply(zeta, I)) \wedge  
 "isafunmode(apply(zeta, I)) \wedge  
 "ispfunmode(apply(zeta, I))  
 by  $saddr(gamma, apply([rho, I], apply2(rho, I), apply3(rho, I)))$ ;

ed08: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkeunop(O, E1)$   
 by  $Rscr(E, zeta, rho, unop(gamma, O))$ ;

ed09: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkebinop(E0, Omega, E1)$   
 by  $Rscr(E0, zeta, rho, Rscr(E1, zeta, rho, binop(gamma, Omega)))$ ;

ed10: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkedere(E)$   
 by  $Rscr(E, zeta, rho, verifyn(gamma))$ ;

ed11: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkefcall(I, Elist)$   
 apply(zeta, I) = mkefunmode  
 by  $Escrstar(Elist, zeta, rho, Sscr(I, gamma))$ ;

ed12: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkefcall(I, Elist)$   
 apply(zeta, I) = mkefunmode(mulist, tau)  
 by  $Aesctr(Elist, zeta, mulist, rho, fcall[apply3(rho, I), gamma])$ ;

ed13: replace  $AEscr(E, zeta, rho, gamma)$   
 where  $E = mkefcall(I, Elist)$   
 apply(zeta, I) = mkefunmode(mulist, tau)  
 by  $Aesctr(Elist, zeta, static_environment, mu, mode)$ ;

## 2. The program

### 2.1. Declarations

pascal	type Continuation	= integer;
	Senvironment	= integer;
	Storage_map	= integer;
	Tcontinuation	= integer;
	Tclist	= (nullist);
	Tenvironment	= integer;
	asyn	= (mkenullist, mkedummy);
	code	= (node);
	compile_environment	= integer;
	function_value	= integer;
	location	= integer;
	mode	= integer;
	static_environment	= integer;
	type	= integer;
	value	= integer;

### 2.2. Virtual procedures

Some of the valuations of the definition of  $L$  Sand  $L7$  are required as virtual functions.

procedure error; entry true; exit false; external :

function Aesctr(E1: asyn; zeta: static\_environment; mu: mode;

```

rho: Senvironment; gamma: Scontinuation; Scontinuation;
entry true; exit true; external ;

function Ccar(G: asyn; zeta: static_environment;
rho: Senvironment; gamma: Scontinuation); Scontinuation;
entry true; exit true; external ;

function Macro(x: code; rho:T;
Tenvironment; gamma:T; Continuation); Continuation;
entry true; exit true; external ;

function Rcar(E0: asyn; zeta: static_environment;
rho: Senvironment; gamma: Scontinuation); Scontinuation;
entry true; exit true; external ;

function Specal(rho: Senvironment; f: asyn;
gamma: Scontinuation); Scontinuation;
entry true; exit true; external ;

function Scarp(f: asyn; gamma: Scontinuation); Scontinuation;
entry true; exit true; external ;

function Binop(gamma: Scontinuation; O: asyn); Scontinuation;
entry true; exit true; external ;

function End(gamma0: gamma); Scontinuation;
entry true; exit true; external ;

function Car(m: mode);
entry true; exit true; external ;

function End(gamma0, gamma); Scontinuation;
entry true; exit true; external ;

function Ctype(mu: mode); ttype;
entry true; exit true; external ;

function D(Dlist: asyn; zeta: static_environment); static_environment;
entry true; exit true; external ;

function Fcall(p: function_value; gamma: Scontinuation); Scontinuation;
entry true; exit true; external ;

function Index(gamma: Scontinuation); Scontinuation;
entry true; exit true; external ;

```

### 2.3. Auxiliary functions

```

function Apply(zeta: static_environment; f: asyn); mode;
entry true; exit true; external ;

function Supply(y: compile_environment; n: integer; f: asyn); location;
entry true; exit true; external ;

function Apply2(rho: Senvironment; f: asyn); integer;
entry true; exit true; external ;

function Apply3(rho: Senvironment; f: asyn); function_value;
entry true; exit true; external ;

```

```

procedure mktAppend(z2: code; s1: code; var z: code;
  rhoT: Tenvironment; gammaT: Tcontinuation);
exit Macr(z, rhoT, gammaT) = Macr(z2, rhoT, Macr(z1, rhoT, gammaT)) ∧
  Lact(z, rhoT, gammaT) =
  append(Lact(z2, rhoT, gammaT), Macr(z1, rhoT, gammaT));
external;

procedure mktLabel(y: compile_environment; l: asyn;
  var ynew: compile_environment);
entry true;
exit disjointLabels(ynew, l, y);
external;

procedure newLabel(var N0: asyn; zeta: static_environment;
  var zeta0: static_environment);
entry true;
exit mkepair(zeta0, N0) = newLabel(zeta);
external;

2.4. Abstract syntax, types and modes
Functions to test, access, and construct objects of abstract syntax, modes
and types are identical to those used in the static semantics. No refinement is
given here.

function isOf/unmode(mu: mode): boolean;
entry true; exit true; external;

function isArrType(tau: itype): boolean;
entry true; exit true; external;

function isAssign(Stmt: asyn): boolean;
entry true; exit true; external;

function isBinOp(E: asyn): boolean;
entry true; exit true; external;

function isVarMode(mu: mode): boolean;

```

```

entry true; exit true; external;

function ifWhile(Stmt: asyn): boolean;
entry true; exit true; external;

function makeBlock(D, G: asyn): asyn;
exit true; external;

function makeCommandList(G0, G1: asyn): asyn;
exit true; external;

...
function mkeEvalMode(tau: itype): mode;
entry true; exit true; external;

procedure matchArrayMode(mu: mode; var muList: mode; var tau: itype);
entry true;
exit mu = mkeArrayMode(muList, tau);
external;

procedure matchVarMode(mu: mode; var tau: itype);
entry true;
exit mu = mkeVarMode(mu, tau);
external;

...
procedure mktArrayType(tau: itype; var nu: integer;
  var tau0: itype; var tau1: itype);
entry true;
exit tau = mkeArrayType(nu, tau0, tau1);
external;

...
procedure matchVarMode(mu: mode; var tau: itype);
exit mkeVarMode(tau) = mu;
external;

procedure matchWhile(Stmt: asyn; var E0: asyn; var G: asyn);
exit mkeWhile(E0, G) = Stmt;
external;

...

```

### 2.5. Code generating functions

We assume the following code generation procedures to be external. An implementation is immediate in most cases.

```

procedure Acode(E: asyn;
               zeta: static_environment;
               mu: mode;
               rho: Senvironment;
               gamma: Scontinuation;
               y: compile_environment;
               rhoT: Tenvironment;
               gammaT: Tcontinuation;
               var z: code);
begin
  entry varthetaU(rho, zeta, yold, rhoT);
  exit varthetaU(rho, y, rhoT)  $\wedge$  rho = Varstar(Dlist, zeta, rho0);

  procedure Bcode(Theta: asyn;
                   zeta: static_environment;
                   rho: Senvironment;
                   gamma: Scontinuation;
                   y: compile_environment;
                   rhoT: Tenvironment;
                   gammaT: Tcontinuation;
                   var z: code);
begin
  initial  z = z0;
  entry varthetaU(rho, zeta, y, rhoT)  $\wedge$ 
        varthetaG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gamma, mat));
  exit varthetaG(Bscr(Theta, zeta, rho, gamma), PhiS(zeta, rho, y),
                Mscr(z, rhoT, gammaT))  $\wedge$ 
      LscrT(z, rhoT, gammaT) = LscrT(z0, rhoT, gammaT);
  external;

  procedure Cicode(T: asyn; zeta: static_environment; rho: Senvironment;
                     gamma: Scontinuation; y: rhoT)  $\wedge$ 
                     rhoT: Tenvironment; gammaT: Tcontinuation; var z: code);
begin
  entry varthetaU(rho, zeta, y, rhoT)  $\wedge$ 
        varthetaG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gamma, mat));
  exit varthetaG(Cscr(T, zeta, rho, gamma), PhiS(zeta, rho, y),
                Mscr(z, rhoT, gammaT));
  external;

  procedure Condcode(N: asyn; zeta: static_environment; rho: Senvironment;
                      gamma: Scontinuation; y: compile_environment;
                      rhoT: Tenvironment; gammaT: Tcontinuation; var z: code);
begin
  entry varthetaU(rho, zeta, y, rhoT)  $\wedge$ 
        varthetaG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gamma, mat));
  exit varthetaG(Cnd(gamma, apply(rho, N)), PhiS(zeta, rho, y),
                Mscr(z, rhoT, gammaT));
  external;

  procedure Dstarcode(Dlist: asyn; zeta: static_environment;
                       rho: Senvironment;
                       gamma: Scontinuation; y: compile_environment;
                       rhoT: Tenvironment; gammaT: Tcontinuation; var z: code);
begin
  initial  rho = rho0;
  entry varthetaU(rho, zeta, yold, rhoT);
  exit varthetaU(rho, y, rhoT)  $\wedge$  rho = Varstar(Dlist, zeta, rho0);
  external;

```

```

entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Dstarstar(Dlist, zeta, rho, gamma), PhiS(zeta, rho, y),
Macr(z, rhoT, gammaT));

external ;

procedure Estartcode(Elist: asyn;
zeta: static_environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile_environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Estarstar(Elist, zeta, rho, gamma),
PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));

procedure Fstartcode(Dlist: asyn; zeta: static_environment;
rho: Senvironment;
y: compile_environment;
rhoT: Tenvironment; gammaT: Tcontinuation; var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(Dstarstar(Dlist, zeta, rho), PhiS(zeta, rho, y), rhoT);
exit varthetaU[rho, zeta, y, rhoT];

external ;

procedure Lcode(E0: asyn; zeta: static_environment;
rho: Senvironment;
gamma: Scontinuation; y: compile_environment;
rhoT: Tenvironment; gammaT: Tcontinuation; var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Lscr(E0, zeta, rho, gamma), PhiS(zeta, rho, y),
Macr(z, rhoT, gammaT));

external ;

procedure Rcode(E: asyn;
G: asyn;
zeta: static_environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile_environment;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

external ;

procedure blockcode(Dlist: asyn;
G: asyn;
zeta: static_environment;
rho: static_environment;
gamma: static_environment;
y: compile_environment;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

external ;

```

```

y: compile_environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code;
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Rscr(E, zeta, rho, gamma),
PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));

procedure Sfcode(I: asyn;
zeta: static_environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile_environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

procedure Spalcode(rho: Senvironment; I: asyn;
gamma: Scontinuation; s: storage_map;
rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

procedure Spalcode(I: asyn; gamma: Scontinuation; s: storage_map;
rho: Senvironment; gammaT: Tcontinuation;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

procedure Sscr(rho: Senvironment; I: asyn;
gamma: Scontinuation; s: storage_map;
rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

procedure blockcode(Dlist: asyn;
G: asyn;
zeta: static_environment;
rho: static_environment;
gamma: static_environment;
y: compile_environment;
var z: code);
entry varthetaU[rho, zeta, y, rhoT] ∧
varthetaG(gamma, PhiS(zeta, rho, y), Macr(z, rhoT, gammaT));
exit varthetaG(Sscr(I, gamma), s, Macr(z, rhoT, gammaT));

external ;

```

```

entry true;
exit forallBLO(Dlist, G, zeta, y, z);
external;

The following functions generate primitive code sequences; i.e. single instructions or instruction sequences to access variables, index arrays and so on.

procedure cmkheader(N: value; rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) = Tapply(Mscr(z0, rhoT, gammaT), epsilon);
external;

procedure cmkeaddr(N: value; rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) = Tunop(Mscr(z0, rhoT, gammaT), O);
external;

procedure cmkebinop(O: asyn; rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) = Taddr(Mscr(z0, rhoT, gammaT), n);
external;

procedure cmkebinop(O: asyn; rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) = Tbinop(Mscr(z0, rhoT, gammaT), O);
external;

procedure cmkeblockdef(znew: code; var z: code;
rhoT: Tenvironment;
gammaT: Tcontinuation);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) =
Mscr[Tmkblock(znew), rhoT, Mscr(z0, rhoT, gammaT)];
external;

procedure cmkelabelcode(n: integer;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) =
LscrT(z, rhoT, gammaT) = Mscr(z, rhoT, gammaT) ∧
cons(Mscr(z0, rhoT, gammaT), LscrT(z0, rhoT, gammaT));
external;

```

```

procedure cmkelti(epsiion: value; rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) = Tapply(Mscr(z0, rhoT, gammaT), epsilon);
external;

procedure cmkeunop(O: asyn; rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit Mscr(z, rhoT, gammaT) = Tunop(Mscr(z0, rhoT, gammaT), O);
external;

procedure fcallcode(l: asyn;
zeta: static-environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile-environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit varthetaG(fcall(apply3(rho, l), gamma), PhiS(zeta, rho, y),
Mscr(z, rhoT, gammaT));
external;

procedure indexcode(zeta: static-environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile-environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit varthetaG(phiG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT)));
external;

procedure indexdecode(zeta: static-environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile-environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
initial z = z0;
entry true;
exit varthetaG(phiG(index(gamma), PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT)));
external;

procedure jumpcode(N: asyn; rho: Senvironment;
gamma: Scontinuation; s: storage-map;
var z: code);
initial z = z0;
entry true;
exit varthetaG(phiG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT)));
external;

procedure jumpcode(N: asyn; rho: Senvironment;
gamma: Scontinuation; s: storage-map;
var z: code);
initial z = z0;
entry true;
exit varthetaG(phiG(index(gamma), PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT)));
external;

```

```

rhoT: Tenvironment; gammaT: Tcontinuation;
zeta: static_environment;
rho: Tenvironment;
gamma: Scontinuation;
y: compile_environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code;

entry varthetaG(gamma, s, Mscr(z, rhoT, gammaT));
exit varthetaG(adjust(apply(rho, N), gamma), s, Mscr(z, rhoT, gammaT));
external;

procedure selectcode(I: asyn;
var z: code);
zeta: static_environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile_environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code;
entry varthetaU(rho, zeta, y, rhoT) ∧
varthetaG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT));
exit varthetaG(select(I, gamma), PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT));
external;

procedure updatecode(gamma: Scontinuation; s: storage-map;
rhoT: Tenvironment; gammaT: Tcontinuation;
var z: code);
entry varthetaG(gamma, s, Mscr(z, rhoT, gammaT));
exit varthetaG(update(gamma), s, Mscr(z, rhoT, gammaT));
external;

procedure verifycode(zeta: static_environment;
rho: Senvironment;
gamma: Scontinuation;
y: compile_environment;
rhoT: Tenvironment;
gammaT: Tcontinuation;
var z: code);
entry varthetaU(rho, zeta, y, rhoT) ∧
varthetaG(gamma, PhiS(zeta, rho, y), Mscr(z, rhoT, gammaT));
exit varthetaG(verify(gamma), y, Mscr(z, rhoT, gammaT));
external;

```

## 2.8. Expressions

```

procedure AEcode,
%has been forwarded%
var E1, E2, Elist, I, N, O: asyn;
procedure AEcode(E: asyn;

```

```

alpha: location;
epsilon: value;
mu,mulist: mode;
nu, m: integer;
tau, tau0, tau1: type;
begin
  if isnumber(E) then
    begin
      matchnumber(E, N);
      cmteit(N, rhoT, gammaT, z);
    end else
      if isezpid(E) then
        begin
          matchezpid(E, I);
          mu ← apply(zeta, I);
          if isfunmode(mu) then
            Sfcode(I, zeta, rho, gamma, y, rhoT, gammaT, z)
          else
            if isunmode(mu) then
              begin
                fcallcode(I, zeta, rho, gamma, y, rhoT, gammaT, z)
              end else
                if isp/unmode(mu) then
                  begin
                    fcallcode(I, zeta, rho, gamma, y, rhoT, gammaT, z)
                  end else
                    if isunmode(mu) then
                      begin
                        m ← apply2(rho, I);
                        alpha ← gossip(y, m, I);
                        cmkeadd(m, rhoT, gammaT, z);
                        cmkelit(alpha, rhoT, gammaT, z);
                      end else
                        if isunop(E) then
                          begin
                            matcharop(E, O, E1);
                            cmkeunop(O, rhoT, gammaT, z);
                            Rcode(E1, zeta, rho, unop(gamma, O), y, rhoT, gammaT, z);
                          end else
                            if isbinop(E) then
                              begin
                                matchbinop(E, E1, O, E2);
                                matchbinop(E, E1, O, E2);
                              end
                            end
                          end
                        end
                      end
                    end
                  end
                end
              end
            end
          end
        end
      end
    end
  end
end

```

```

cmkebinop(O, rhoT, gammaT, z);
Rcode(E2, zeta, rho, binop(gamma, O), y, rhoT, gammaT, z);
Rcode(E1, zeta, rho, Rcr(E2, zeta, rho, binop(gamma, O)),
y, rhoT, gammaT, z);
end else
if isderf(E) then
begin
  matchderf(E, E1);
  verifycode(zeta, rho, gamma, y, rhoT, gammaT, z);
  Rcode(E1, zeta, rho, verify(gamma), y, rhoT, gammaT, z);
end else
if isfcall(E) then
begin
  matchfcall(E, I, Elist);
  matchfcall(E, I, Elist);
  mu ← apply(zeta, I);
  if isfunmode(mu) then
    begin
      Sfcode(I, zeta, rho, gamma, y, rhoT, gammaT, z);
      Ecoderar(Elist, zeta, rho,
Sactf(I, gamma), y, rhoT, gammaT, z);
    end else
      if isp/unmode(mu) then
        begin
          matchafunmode(mu, multet, tau);
          fcallcode(I, zeta, rho, gamma, y, rhoT, gammaT, z);
          Acoderar(Elist, zeta, multet, rho,
Sactf(I, gamma), y, rhoT, gammaT, z);
        end
      end
    end
  begin
    if isfunmode(mu) then
      begin
        matchafunmode(mu, multet, tau);
        fcallcode(I, zeta, rho, gamma, y, rhoT, gammaT, z);
        Acoderar(Elist, zeta, multet, rho,
fcall(apply3(rho, I), gamma),
y, rhoT, gammaT, z);
      end else
        if isp/unmode(mu) then
          begin
            matchafunmode(mu, multet, tau);
            fcallcode(I, zeta, rho, gamma, y, rhoT, gammaT, z);
            Acoderar(Elist, zeta, multet, rho,
fcall(apply3(rho, I), gamma),
y, rhoT, gammaT, z);
          end
        end
      end
    end
  begin
    if isfcall(E, E1) then
      begin
        matchfcall(E, E1, E2);
        mu ← cc(E1, zeta);
      end
    end
  end
end

```

```

tau ← ctyp{mu};
if isarraytype(tau) then
begin
matcharraytype(tau, nu, tau0, tau1);
indecod(zeta, rho, gamma, y, rhoT, gammat, z);
Acode(E2, zeta, mkevalmode(tau0), rho,
      indez(gamma), y, rhoT, gammat, z);
Bcode(E1, zeta, rho,
      Acr(E2, zeta, mkevalmode(tau0), rho,
           indez(gamma)),
      y, rhoT, gammat, z);
end else error
end else
if isselect(e) then
begin
matchselect(E, E1, I);
selectedcode(I, zeta, rho, gamma, y, rhoT, gammat, z);
Bcode(E1, zeta, rho, select(I, gamma), y, rhoT, gammat, z);
end else error;
end;
end;

```

**2.7. Commands**

```

procedure Ccode(G: asyn;
                zeta: static_environment;
                rho: Senvironment;
                gamma: Scontinuation;
                y: compile_environment;
                rhoT: Tenvironment;
                gammat: Tcontinuation;
                var z: code);
initial z = z0;
entry varthetaU(rho, zeta, y, rhoT) ∧
varthetaG(gamma, PhiS(zeta, rho, y), gammat);
exit  varthetaG(Cacr(G, zeta, rho, gamma), PhiS(zeta, rho, y),
               Mscr(z, rhoT, gammat)) ∧
varthetaGstar(Jscr(G, zeta, rho, gamma), PhiS(zeta, rho, y),
              LscrT(z, rhoT, gammat));
var G0, G1, Theta, N: asyn;
m: integer;
x1, x2: code;
begin
if islabel(G) then
begin
mdchlabel(G, N, Theta);
z ← nocode;
Bcode(Theta, zeta, rho, gamma, y, rhoT, gammat, z);
m ← apply(y, n);
cmklabelcode(m, rhoT, gammat, z);
end else
if isstmt(G) then
begin
matchstmt(G, Theta);
z ← nocode;
Bcode(Theta, zeta, rho, gamma, y, rhoT, gammat, z);
end else
if iscommandlist(G) then
begin
matchcommandlist(G, GU, GT);

```

```

Code(G1, zeta, rho, gamma, y, rhoT, gammat, z);
Code(G0, zeta, rho, Cerr(G1, zeta, rho, gamma), y, rhoT,
    Macr(z, rhoT, gammaT), s2);
makeTappend(s2, s1, z, rhoT, gammaT);
end else error
end;

procedure Gcode(G: asyn;
    zeta: static_environment;
    y: compile_environment;
    var z: code);
begin
initial z = z0;
exit forallUGG(G, zeta, y, s);
external ;

```

**2.8. Statements**

```

procedure Bcode(Stmt: asyn;
    zeta: static_environment;
    rho: Senvironment;
    gamma: Scontinuation;
    y: compile_environment;
    rhoT: Tenvironment;
    gammat: Tcontinuation;
    var z: code);
entry varthetaU[rho, zeta, rho, rhoT] ∧
varthetaG[gamma, PhiS[zeta, rho, y], Macr(z, rhoT, gammat)];
exit varthetaG[Back(Stmt, zeta, rho, gamma), PhiS[zeta, rho, y],
    Macr(z, rhoT, gammaT)];
var Dist, E0, E1, Elst, G, G0, G1, I, N, N0, N1, T: asyn;
m: integer;
mu, mult: mode;
tau: ttype;
zeta0, zeta1: static_environment;
znew: code;
ynew: compile_environment;
begin
if isassign(Stmt) then
begin
matchassign(Stmt, E0, E1);

```

```

    mu ← ce(E0, zeta);
    if isuarmode(mu) then
begin
matchuarmode(mu, tau);
updatecode(mu, PphiS(zeta, rho, y), rhoT, gammat, z);
Acode(E1, zeta, mkeuarmode(tau), rho, update(gamma), y, rhoT, gammat, z);
Lcode(E0, zeta, rho, Ascr(B), zeta, mkewalmodetau, rho,
    y, rhoT, gammat, z);
update(gamma));
end else error
end
else
if isunless(Stmt) then
begin
matchunless(Stmt, E0, N);
Concode(N, zeta, rho, gamma, y, rhoT, gammat, z);
Rcode(E0, zeta, rho, Cnd(gamma, apply(rho, N)), y, rhoT, gammat, z);
end
else
if iss(Stmt) then
begin
matchiss(Stmt, E0, G0, G1);
newsourcelabel(N0, zeta, zeta0);
newsourcelabel(N1, zeta0, zeta1);
T ← mkecommandist(mkestmt(mkeunless(E0, N0)),
mkecommandist(G0,
mkecommandist(mkestmt(mkegoto(N1)),
mkecommandist(mklabel(N0, mkstmt(mkblock(mkenulllist, G1)))),
mkestmt(mtedummy))))));
Ccode(T, zeta, rho, gamma, y, rhoT, gammat, z);
end
else
if isdum,ny(Stmt) then begin end else
if isuwhile(Stmt) then
begin
matchuwhile(Stmt, E0, G);
newsourcelabel(N0, zeta, zeta0);
newsourcelabel(N1, zeta0, zeta1);
T ← mkecommandist(mkelabel(N0, mkunless(E0, N1)),
mkecommandist(G,
mkecommandist(mkstmt(mkegoto(N0)),
mkelabel(N1, mkedummy))))));
Ccode(T, zeta, rho, gamma, y, rhoT, gammat, z);
end
else
if isassign(Stmt) then
begin
matchassign(Stmt, E0, E1);

```

```

if isrepeat(Stmt) then
begin
matchrepeat(Stmt, G, EO);
newsourcelabel(N0, zeta, zeta0);
T ← makecommandlist(makelabel([N0, makeblock(mkenullist, G)]);
mkeitml(mkeunless([EO, N0]));
Gicode[T, zeta0, rho, gamma, y, rhoT, gammat, z];
end else
if isgoto(Stmt) then
begin
matchgoto(Stmt, N);
jmpcode[N, rho, gamma, PhiS(zeta, rho, y), rhoT, gammat, z];
end else
if ispcall(Stmt) then
begin
matchpcall(Stmt, I, Elist);
mu ← apply(zeta, I);
if isprecond(mu) then
begin
Scrpode[I, gamma, PhiS(zeta, rho, y), rhoT, gammat, z];
Estarcod[Elist, zeta, rho, Scrp[I, gamma],
y, rhoT, gammat, z];
end else
begin
cestar(Elist, zeta, muist, m);
Scrpode[rho, I, gamma, PhiS(zeta, rho, y), rhoT, gammat, z];
Acoderar[Elist, zeta, muist, rho, Scrp[I, gamma],
y, rhoT, gammat, z];
end
end else
if isblock(Stmt) then
begin
matchblock(Stmt, Dlist, G);
Allocate[Dlist, zeta, rhoT, rho, y, ynew];
zeta0 ← d(Dlist, redeftrue[zeta, j(G)]);
blockcode[Dlist, G, zeta0, ynew, znew];
cmkblockcode[snew, z, rhoT, gammat];
end else error;
end ;

```